

Fundamentals of Dynamical Systems (Tópicos de Sistemas Dinâmicos) Licenciatura em Matemática

Salvatore Cosentino

Departamento de Matemática - Universidade do Minho

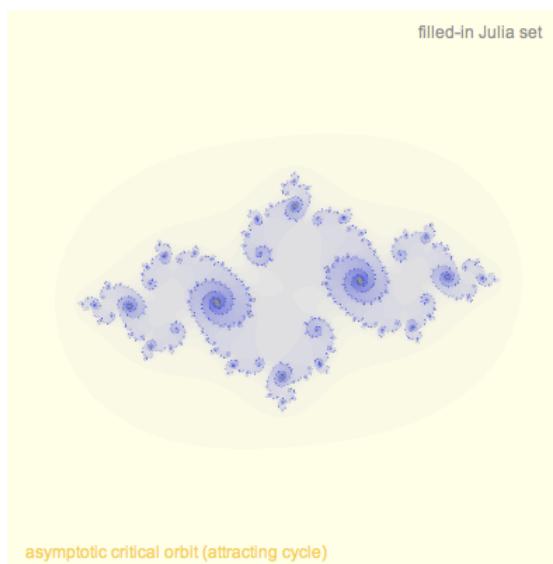
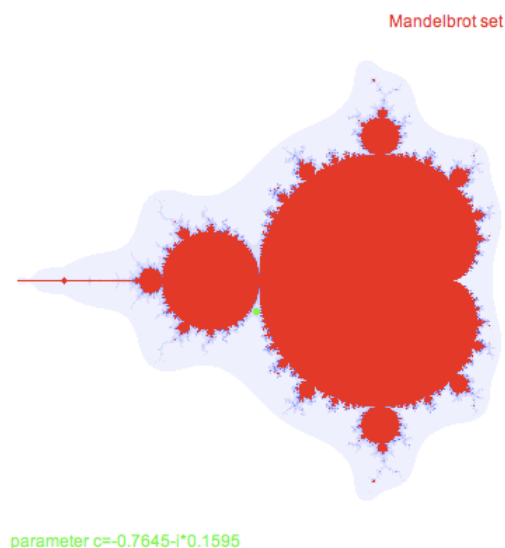
Campus de Gualtar - 4710 Braga - PORTUGAL

gab: CG - Edifício 6 - 3.48, tel: 253 604086

e-mail: scosentino@math.uminho.pt

url <http://w3.math.uminho.pt/~scosentino>

May 18, 2021



This work is licensed under a
[Creative Commons Attribution-ShareAlike 3.0 Unported License](https://creativecommons.org/licenses/by-sa/3.0/).

Abstract

This is an english version of the notes written for my lectures on “Tópicos de Sistemas Dinâmicos” for the “Licenciatura em Matemática” of the University of Minho, during the last decade (available at my page <http://w3.math.uminho.pt/~scosentino/salteaching.html>). Emphasis is on examples presenting in the simplest way some important ideas, and on the interplay between different areas of mathematics. Some very important parts of the modern theory of dynamical systems, as hyperbolic theory, Lyapunov exponents, Hamiltonian systems, thermodynamic formalism, ... are almost completely missing. Other interesting results or directions are only sketched.

Classical modern references and sources are the encyclopedic [KH95] and the introductory [HK03], by Anatole Katok and Boris Hasselblat,. I also recommend the wonderful set of notes [Kn05] by Oliver Knill. Others references are suggested along the text.

It would be nice to have time and places to do simulations, using some of the software at our disposal in laboratories: this includes proprietary software, like [Mathematica®8](#) and [Matlab](#), or open software, like [Python](#) and [GeoGebra](#). Occasionally, we may also use some [c++](#) code and [Java](#) applets. Some applets are in the [bestiario](#) in my [web page](#), and everything about the course may be found in my pages

<http://w3.math.uminho.pt/~scosentino/salteaching.html>

Black paragraphs form the main text.

Blue paragraphs are important or interesting examples, or computations, most of them even more important than black paragraphs.

Red paragraphs are non-trivial facts and results which may be skipped in a first (and also second) reading.

ex: means “exercise”, to be solved at home or in the classroom.

A \square indicates the end of a proof.

Pictures were made with *Grapher* or [SketchBook](#) on my MacBook, or taken from [Wikipedia](#), or produced with [Matlab](#), [Python](#) or [Java](#) codes, like the one in the front page.

Contents

1	Iterations	5
1.1	Exponential growth/decay	5
1.2	Linear recursions	11
1.3	Iteration of maps	13
1.4	Babylonians-Heron method to compute square roots	17
1.5	From Newton method to Julia and Fatou sets	20
2	Differential equations and flows	23
2.1	Flows	23
2.2	Structure of physical models	25
2.3	Integration of one-dimensional systems	27
2.4	Existence and uniqueness theorems	31
2.5	Oscillations and cycles	35
2.6	Phenomenological models	38
3	Topological dynamical systems, basic definitions	43
3.1	Transformations	43
3.2	Trajectories and orbits	44
3.3	Periodic orbits and basin of attraction	45
3.4	Observables	49
3.5	Invariant sets	50
3.6	Conjugations	50
4	Linear systems	52
4.1	Exponential of a linear operator	52
4.2	Linear flows	54
4.3	Linear systems in the plane	56
4.4	Jordan normal form	59
4.5	Hyperbolic linear flows	63
5	Numbers and dynamics	65
5.1	Decimal expansion and multiplication by ten	65
5.2	Bernoulli shifts	68
5.3	Rotations of the torus	69
5.4	Dyadic adding machine	71
5.5	Continued fractions and Gauss map	72
5.6	Exponential sums	76
6	Simple orbits and perturbations	78
6.1	Topological fixed point theorems	78
6.2	Dynamics of contractions	78
6.3	Linear maps	81
6.4	Order of the line and trajectories	84
6.5	Local analysis: attracting and repelling fixed points	85
6.6	Transversality and bifurcations	89
7	Statistical description of orbits	92
7.1	Probability measures	92
7.2	Transformations and invariant measures	94
7.3	Invariant measures and time averages	97
7.4	Examples of invariant measures	99

8	Recurrences	102
8.1	Limit sets and recurrent points	102
8.2	Dirichlet theorem on Diophantine approximation	103
8.3	Poincaré recurrence theorem	104
8.4	Transitivity and minimality	106
8.5	Kronecker theorem on irrational rotations	108
8.6	Circle homeomorphisms	110
9	Chaos	114
9.1	Sensitive dependence on initial conditions	114
9.2	Topological mixing	116
9.3	Expanding maps of the circle	118
9.4	Symbolic dynamics and coding	120
9.5	Non-negative matrices and the Perron-Frobenius theorem	125
9.6	Cantor sets	131
9.7	Hyperbolic automorphisms of the torus	133
9.8	Horseshoes and solenoids	135
10	Topological entropy and zeta function	139
10.1	Topological entropy	139
10.2	Expansiveness and generators	140
10.3	Dimensions of metric spaces	143
10.4	Topological entropy according to Bowen and Dinaburg	144
10.5	Growth of periodic orbits and zeta function	148
11	Ergodicity and convergence of time means	152
11.1	Ergodicity	152
11.2	Examples of ergodic maps	154
11.3	Normal numbers	157
11.4	Distribution of digits in continued fractions	158
11.5	Unique ergodicity and equidistribution	160
11.6	Mixing	162

1 Iterations

1.1 Exponential growth/decay

Fibonacci numbers. Consider the following problem, posed by Leonardo Pisano (alias Fibonacci) in his *Liber Abaci*, 1202:

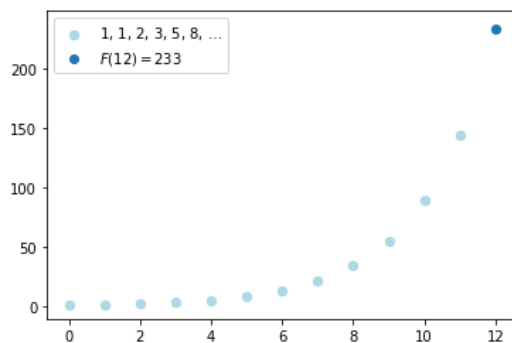
*Quot paria cuniculorum in uno anno ex uno pario germinentur.
Quidam posuit unum par cuniculorum in quodam loco, qui erat undique pariete circumdatus, ut sciret, quot ex eo paria germinarentur in uno anno: cum natura eorum sit per singulum mensem aliud par germinare; et in secundo mense ab eorum nativitate germinant.*

Let f_n be the number of pairs of rabbits at the n -th month. The offspring one month later, $f_{n+1} - f_n$, is equal to the number of “adult” pairs present in the n -th month, which is f_{n-1} . Therefore, the f_n ’s satisfy the recursive law

$$f_{n+1} = f_n + f_{n-1}, \quad (1.1)$$

which prescribes the successive values of f_n given some initial values f_0 and f_1 . The sequence grows quite fast, as you can see: if we take, with Fibonacci, the initial values $f_0 = f_1 = 1$, we get

1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, 233, 377, 610, 987, 1597, 2584, ...

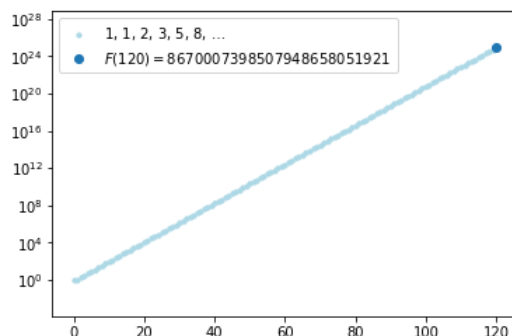


First 13 Fibonacci numbers.

These numbers soon become astronomically large. For example, after 10 years we get

$$f_{120} \simeq 8.67 \times 10^{24},$$

larger than the Avogadro number! In order to see their growth, we must use a logarithmic scale.



First 121 Fibonacci numbers in logarithmic scale.

An [applet](http://w3.math.uminho.pt/~scosentino/salbestiario.html) which computes the sequence is in my page <http://w3.math.uminho.pt/~scosentino/salbestiario.html>. Also useful would be a formula, or at least an asymptotic formula, for the f_n ’s, and I’ll show you one later. For example, an asymptotic formula would solve a problem like

ex: Estimate the smallest time n such that $f_n > 10^{80}$.

Duplication of bacteria. Experiments show that a population of bacteria, during a certain initial period at least, double each characteristic time $\tau > 0$. Thus, an initial population of N_0 cells gives origin to $N_1 = 2N_0$ after a time τ , to $N_2 = 4N_0$ cells after a time 2τ , ..., and to

$$N_n = 2^n N_0$$

cells after time $n\tau$. For example, a unique cell gives origin to 1024 cells after a time $t = n\tau$ such that $2^n = 1024$, i.e. $n\tau = (\log_2 1024) \cdot \tau = 10 \cdot \tau$.

Sequences as time series. A (real or complex valued) *sequence* is a collection $(x_n)_{n \in \mathbb{N}_0}$ of numbers $x_n \in \mathbb{R}$ or \mathbb{C} , indexed (hence ordered) by a non-negative integer $n \in \mathbb{N}_0 := \{0, 1, 2, 3, \dots\}$. We may think of the index n as “time”, and therefore at the n -th term x_n as the value of some “observable” x at time n (as the number of pairs of rabbits or of bacteria). Physicists call them “time series”.

Clearly, we may as well define sequences with values in an arbitrary set X , for example in the Euclidean space \mathbb{R}^d . Also, we may allow time n to be negative, for example to live in \mathbb{Z} . Such collections $(x_n)_{n \in \mathbb{Z}}$ are called “two-sided” sequences.

Subsequences are obtained forgetting to observe x at certain times, i.e. are sequences $(y_i)_{i \in \mathbb{N}_0}$ defined by $y_i := x_{n_i}$, where $i \mapsto n_i$ is an increasing function of \mathbb{N}_0 into itself.

Sequences may be defined as functions are. Indeed, a sequence with values in the set X is nothing but a function

$$x : \mathbb{N}_0 \rightarrow X,$$

disguised by the notation $x_n := f(n)$. A second possibility, more interesting for our point of view, is some recursive law prescribing the value of x_n given the (past) values of x_0, x_1, \dots, x_{n-1} . A third possibility, is using some property that the successive terms must have.

Engineers also use to look at sequences as “discrete-time signals” $x[n] = x(n\tau)$, possibly obtained from an analogic signal $x(t)$, defined for times t in some interval of the real line, sampling its values at integer multiples of some “sampling time” τ .

Discrete derivative and primitives. Given a sequence $x = (x_n)$ with values, for example, in \mathbb{R} , we can “integrate” and get the sequence Sx , defined by $(Sx)_0 := 0$ and

$$(Sx)_n := \sum_{k=0}^{n-1} x_k = x_0 + x_1 + \dots + x_{n-1} \quad \text{for } n \geq 1.$$

The operator S should be called *sum operator*, or also *discrete primitive*. Also, we may define its (*forward*) *discrete derivative* taking differences, as the sequence Dx defined by

$$(Dx)_n := x_{n+1} - x_n.$$

It is clear that S and D are discrete versions of the integration and derivation operators, respectively. Indeed, one easily checks that Newton’s fundamental theorem of calculus and Leibniz rule looks like

$$(DSx)_n = x_n \quad \text{and} \quad (SDx)_n = x_n - x_0,$$

respectively (observe that D and S do not commute, and that their commutator $[D, S]$ is the operator sending a sequence x to the constant sequence equal to x_0). A law like $(Sx)_n = b_n$ should be thought as a discrete “differential equation” solvable for the first derivative, and it is indeed solved by integration, as $x_n = x_0 + (Sb)_n$.

Arithmetic progression. An *arithmetic progression* $x_n = a + nb$, which may also be defined using the recursion $x_{n+1} = x_n + b$, with initial term $x_0 = a$. It is a solution of the “discrete differential equation” $Dx = b$ (the constant sequence $b_n = b$) with initial condition $x_0 = a$.

ex: Consider the sequence 1, 2, 5, 10, 17, 26, ... Differentiate twice, and guess the next term.

ex: Consider the sequence 1, 2, 9, 28, 65, 125, ... Differentiate enough times to guess the next term.

Discrete exponential. A discrete version of the differential equation $\dot{x} = x$ satisfied by the exponential function $x(t) = e^t$ is $Dx = x$. Can you recognise the sequence it generate?

Fibonacci difference equation. The Fibonacci sequence $1, 1, 2, 3, 5, 8, \dots$ satisfies

$$(Df)_n = f_{n-1}$$

(where we set $f_{-1} := 0$), a shifted version of the equation $Dx = x$.

The primes sequence. The sequence $2, 3, 5, 7, 11, 13, 17, 19, 23, \dots$, whose generic term is the n -th prime number p_n . It is not clear what the recursive law could be.

Euler method and discrete differential equations. Suppose we have a autonomous differential equation

$$\dot{y} = v(y)$$

defined by a vector field $v(y)$. Following Euler, we may discretize time looking at the variable at integer multiples $t = n\tau$ of some fixed (possibly small but positive) “time-step” τ . The solution $y(n\tau)$ is then approximated with the sequence x_n , defined by the recursion

$$x_{n+1} - x_n \simeq \tau \dot{y}(n\tau) = \tau v(y(n\tau)) \simeq \tau v(x_n)$$

provided some initial condition $x_0 = y(0)$. The above is a discrete differential equation of the form $(Dx)_n = F(x_n)$.

Limits. We say that the real or complex sequence (x_n) *converges* to some *limit* $a \in \mathbb{R}$ or \mathbb{C} , and we write $\lim_{n \rightarrow \infty} x_n = a$ or simply $x_n \rightarrow a$ (as $n \rightarrow \infty$), if for any “precision” $\varepsilon > 0$ there exists a time \bar{n} such that $|x_n - a| < \varepsilon$ for all times $n \geq \bar{n}$. This means that the values x_n are within an arbitrarily small neighbourhood of a as long as the time n is sufficiently large.

The basic fact about limits in the real line \mathbb{R} is that monotone (non-decreasing or non-increasing, i.e. satisfying $x_{n+1} \geq x_n$ or $x_{n+1} \leq x_n$, for any n , respectively) bounded (i.e. such that $|x_n| \leq M$ for some $M > 0$ and all n) sequences of real numbers do admit limit. For example, the limit of a bounded increasing sequence is simply the supremum of the set of values.

We also use the notation $x_n \rightarrow \pm\infty$ to say that given an arbitrarily large $K > 0$ we can find a time \bar{n} such that $\pm x_n > K$ for all times $n \geq \bar{n}$.

Of course, there exist sequences which do not admit limits in either senses. These are, for example, oscillating sequences, as $x_n = (-1)^n$. We’ll encounter sequences with stranger behaviors.

Fundamental sequences. A sequence (x_n) is said *fundamental*, or *Cauchy sequence*, if for any precision $\varepsilon > 0$ there exists a time \bar{n} such that

$$|x_n - x_m| < \varepsilon$$

for all times $n, m > \bar{n}$. Fundamental sequences are clearly bounded. It is obvious that a convergent sequence is fundamental (a triangular argument, since both x_n and x_m are $\varepsilon/2$ -near to the limit for sufficiently large n and m). A similar triangular argument shows that a fundamental sequence with a convergent subsequence is itself convergent. Less obvious is that any fundamental sequence in \mathbb{R} is convergent. Indeed, let $X_n := \{x_k \text{ with } k \geq n\}$. It is clear that the X_n are bounded, and therefore by the supremum axiom there exist the numbers $a_n := \inf X_n$. But the sequence (a_n) is bounded and not decreasing, and therefore there exists $a = \lim_{n \rightarrow \infty} a_n$ (indeed, $a = \sup \{a_n \text{ with } n \in \mathbb{N}\}$). It is then easy to construct subsequences of (x_n) which converge to a , and this implies that (x_n) itself is convergent to a .

Thus, we may know that a sequence is convergent without knowing its limit! In general, convergence of all fundamental sequences is taken as a definition of (sequential) completeness of a metric space.

Geometric progression. The most important sequence is the *geometric progression*, defined by the recursion

$$x_{n+1} = \lambda x_n,$$

and an initial term $x_0 = a$ (which we may assume $\neq 0$ to avoid trivialities). Thus, the sequence is

$$x_0 = a \quad x_1 = a\lambda \quad x_2 = a\lambda^2 \quad \dots \quad x_n = a\lambda^n \quad \dots$$

The parameter λ (which may be real or complex) is called *ratio*, since it is the ratio x_{n+1}/x_n between successive terms of the sequence.

The geometric sequence clearly converges to zero when $|\lambda| < 1$. It is constant, hence trivially convergent, when $\lambda = 1$, while oscillates between $\pm a$ when $\lambda = -1$ (hence does not converge if $a \neq 0$). We may also observe that $|\lambda^n| \rightarrow \infty$ when $|\lambda| > 1$.

ex: Show that the geometric progression $x_n = a\lambda^n$ is the solution of the discrete autonomous differential equation $Dx = \gamma x$ (doesn't it remind the differential equation defining the exponential?) with initial condition $x_0 = a$, where the parameter is $\gamma = \lambda - 1$. In particular, verify that the doubling progression $x_n = 2^n$ satisfies $Dx = x$ with initial condition $x_0 = 1$.

ex: Show that each term $x_n = a\lambda^n$ of a geometric progression is equal to the geometric mean $\sqrt{x_{n+1}x_{n-1}}$ of its neighbors (provided $n > 0$, of course).

Computing limits. First, observe that $x_n \rightarrow a$ is equivalent to $x_n - a \rightarrow 0$. Therefore, we only need to understand how to “prove” that some sequence converges to zero. One possibility is to “compare” the sequence (x_n) under investigation with a sequence with known behavior, as for example the geometric progression. Indeed, if $|x_n| \leq y_n$ for all n sufficiently large, then $y_n \rightarrow 0$ implies $x_n \rightarrow 0$ too.

Subsequences and sequential compactness. A *subsequence* of a sequence (x_n) is a sequence (x_{n_i}) obtained selecting only the values x_{n_i} of the original sequence, where $i \mapsto n_i$ is an increasing map $\mathbb{N}_0 \rightarrow \mathbb{N}_0$.

Sometimes we are only interested in a rough estimate of the growth of a sequence (x_n) . The “limsup” is the limit $\limsup_{n \rightarrow \infty} x_n := \lim_{n \rightarrow \infty} a_n \in \mathbb{R} \cup \{\infty\}$ of the non-increasing sequence $a_n := \sup\{x_n, x_{n+1}, x_{n+2}, \dots\}$. The “liminf” is the limit $\liminf_{n \rightarrow \infty} x_n := \lim_{n \rightarrow \infty} b_n \in \mathbb{R} \cup \{-\infty\}$ of the non-decreasing sequence $b_n := \inf\{x_n, x_{n+1}, x_{n+2}, \dots\}$.

The basic fact (that closed and bounded sets of the real line are *sequentially compact*) is that any bounded sequence admits a convergent subsequence.

Exponential decay and half-life. Radioactive decay may be characterized by a “half-life”, the time τ after which approximately half of the initial nuclei decay, between a sufficiently large sample. If q_n denotes the number of nuclei at time $n\tau$, with $n = 0, 1, 2, \dots$, then

$$q_{n+1} = \frac{1}{2} q_n.$$

Thus, the number of nuclei at time $n\tau$ is $q_n = 2^{-n} q_0$, while the product of the decay is $q_0 - q_n = q_0(1 - 2^{-n})$. Observe that $q_n \rightarrow 0$ when $n \rightarrow \infty$.

If solar radiation produces radioactive nuclei at a constant rate $\alpha > 0$ (i.e. α nuclei each time interval τ), then the number of radioactive nuclei at time $n\tau$ satisfies the recursion

$$q_{n+1} = \frac{1}{2} q_n + \alpha. \quad (1.2)$$

Equilibrium is possible when q_0 is equal to $\bar{q} := 2\alpha$, since then $q_1 = \alpha + \alpha = q_0$, $q_2 = \alpha + \alpha = q_1 = q_0$, \dots and so on, $q_n = \bar{q}$ for all $n \in \mathbb{N}$.

What happens if the initial condition is $q_0 \neq \bar{q}$? The recursion says that

$$\begin{aligned} q_1 &= \frac{1}{2}q_0 + \alpha \\ q_2 &= \frac{1}{4}q_0 + \frac{1}{2}\alpha + \alpha \\ q_3 &= \frac{1}{8}q_0 + \frac{1}{4}\alpha + \frac{1}{2}\alpha + \alpha \\ &\vdots \\ q_n &= \frac{1}{2^n}q_0 + \left(\frac{1}{2^{n-1}} + \cdots + \frac{1}{8} + \frac{1}{4} + \frac{1}{2} + 1\right)\alpha \end{aligned}$$

The first term $2^{-n}q_0 \rightarrow 0$ when $n \rightarrow \infty$, which means that “future” is independent on the initial condition q_0 . The second term converges to the equilibrium $\bar{q} = 2\alpha$ when $n \rightarrow \infty$, the factor of α being the sum of a geometric series of ratio $1/2$ (if you forgot about it, see below).

A simpler formula, and insight, for q_n may be obtained using the substitution $x_n := q_n - \bar{q}$, where $\bar{q} = 2\alpha$ is the equilibrium solution. We get

$$\begin{aligned} x_{n+1} &= q_{n+1} - 2\alpha \\ &= \frac{1}{2}q_n + \alpha - 2\alpha \quad (\text{using (1.2)}) \\ &= \frac{1}{2}x_n, \end{aligned}$$

So, the difference between q_n and \bar{q} is a geometric progression with ratio $1/2$. Thus $x_n = x_0 2^{-n}$, and therefore

$$q_n = 2\alpha + (q_0 - 2\alpha) \cdot 2^{-n}.$$

Again, it is interesting to observe that $x_n \rightarrow 0$, and therefore $q_n \rightarrow \bar{q}$, when $n \rightarrow \infty$. So, the amount of radioactive nuclei converges to the stationary value independently on its initial value.

ex: After how much time does the radioactive substance decrease to $\frac{1}{32}$ -th of its initial value?

ex: Half-life of ^{14}C is estimated to be $\tau \simeq 5730$ years. Show how to date a fossil, assuming that we know the proportion of ^{14}C in a living being. ¹

Exponential growth. Exponential growth of populations in a illimited environment is modeled by the recursion

$$p_{n+1} = \lambda p_n,$$

where p_n represent the population at time n (measured in units of some fixed time interval $\tau > 0$), given an initial population p_0 . The meaning of the parameter λ is the following: at every time interval τ , the increase $p_{n+1} - p_n$ of the population is equal the “offspring” αp_n , where $\alpha > 0$ is some “fertility” coefficient, minus the “deaths” βp_n , where $\beta > 0$ is some “mortality” coefficient. Thus, $\lambda = \alpha - \beta$. An applet with the simulations is in [exponentialgrowth](#).

ex: Discuss the behaviour of solutions p_n for different values of λ .

ex: To a population growing exponentially is added or retired a certain amount β each time interval τ . Te model is therefore

$$p_{n+1} = \lambda p_n + \beta,$$

where β is a positive or negative parameter. Find the stationary solution, and then the solution with arbitrary initial condition p_0 (consider the substitution $x_n = p_n - \bar{p}$, where \bar{p} is the stationary solution).

ex: For which values of λ and β do the solution p_n converge to the stationary solution when $n \rightarrow \infty$?

¹J.R. Arnold and W.F. Libby, Age determinations by Radiocarbon Content: Checks with Samples of Known Ages, *Sciences* **110** (1949), 1127-1151.

Growth of Fibonacci numbers. How fast do Fibonacci numbers grow? Define the quotients $q_n := f_{n+1}/f_n$ between neighbor Fibonacci numbers. From (1.1) one deduce the recursive equation

$$q_{n+1} = 1 + 1/q_n \quad (1.3)$$

for the q_n 's. We compute:

$$1, \quad 2, \quad 3/2 = 1.5, \quad 5/3 \simeq 1.66666, \quad 8/5 = 1.6, \quad 13/8 = 1.625, \quad 21/13 \simeq 1.61538, \quad \dots$$

You may observe the sequence in the following [applet](#). It turns out that the sequence (q_n) converge (try to prove it!), namely, $q_n \rightarrow \phi$ as $n \rightarrow \infty$. Taking the limits in the recursive equation (1.3) we see that $\phi = 1 + 1/\phi$, and therefore ϕ is the positive root of the quadratic polynomial $x^2 - x - 1$,

$$\phi = \frac{1 + \sqrt{5}}{2} \simeq 1.6180339887498948482 \dots$$

Hence, for large values of n we may approximate Fibonacci law as

$$f_{n+1} \approx \phi f_n,$$

an exponential growth with rate ϕ . In particular, we expect $f_n \sim \phi^n$.

The limit ϕ is a famous irrational, the Greeks' "[ratio/proportion](#)". As described by Euclid²:

"A straight line is said to have been cut in extreme and mean ratio when, as the whole line is to the greater segment, so is the greater to the less."

If a is the greater part and b the less of a line of length $a + b$, Euclid's requirement is

$$\frac{a+b}{a} = \frac{a}{b}$$

There follows that the ratio $\phi = a/b$ satisfies $1 + 1/\phi = \phi$. This division of an interval is used in Book IV of the Elements to construct a regular pentagon. Observe that, as follows from the quadratic equation, ϕ^{-1} is equal to $\phi - 1$.

ex: Show that ϕ is irrational using its geometric definition (see Euclid's *Elements*, or [\[HW59\]](#) section 4.6.)

The invention of chess. Legend says that Sissa invented chess, and offered the game to the king of Persia. Asked for a reward, he suggested that he wanted one grain of rice on the first square of the chessboard, two grains on the second, four grains on the third, and so on. The king didn't take it seriously, but a computation shows that the reward amounts to

$$1 + 2 + 4 + 8 + \dots + 2^{63} \simeq 1.84 \times 10^{19}$$

grains of rice. Now, if 1 Kg of rice contains something like 30000 grains, the above number amounts to roughly 6.13×10^{11} tons of rise (which you may want to compare with People's Republic of China's production in 2017, which has been, according to [FAO](#), about 2.14×10^8 tons!).

Series. A *series* is a formal infinite sum $\sum_{n=0}^{\infty} x_n$, or $\sum_{n \geq 0} x_n$, where the $x_n \in \mathbb{R}$ are elements of some given real (or complex) sequence. If the sequence (s_n) of *partial sums*, defined as $s_n := \sum_{k=0}^n x_k$ (which are honest numbers) converges to some limit, say $\lim_{n \rightarrow \infty} s_n = s$, then we say the series is *convergent* (or *summable*), and that its *sum* is $\sum_{n \geq 0} x_n := s$.

A series $\sum_n x_n$ is *absolutely convergent* is the series $\sum_n |x_n|$, formed with the absolute values of its terms, is convergent. Of course, absolute convergence is stronger than mere convergence. Indeed, convergent but not absolutely convergent series are quite interesting and strange objects³ (see, for example, the last book by Hardy [\[Har49\]](#)).

²Euclid, *Elements*, Book VI, Definition 3.

³According to Abel (1828), "divergent series are the invention of the devil, and it is shameful to base on them any demonstration whatsoever."

Harmonic series. The *harmonic series*

$$\sum_{n=1}^{\infty} \frac{1}{n} = 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \dots$$

diverges. Indeed, its generic term $1/n$, for $n \geq 1$, is bigger than the integral $\int_n^{n+1} dx/x$, hence the partial sums $\sum_{k=1}^n 1/k$ are greater than the logarithm $\log(n+1)$.

Geometric series. The identity $(1 + \lambda + \lambda^2 + \lambda^3 + \dots + \lambda^n)(\lambda - 1) = \lambda^{n+1} - 1$ shows that, if $\lambda \neq 1$, the sum of the first $n+1$ terms of the geometric progression (with $a = 1$) is

$$1 + \lambda + \lambda^2 + \lambda^3 + \dots + \lambda^n = \frac{\lambda^{n+1} - 1}{\lambda - 1}$$

In particular, when $|\lambda| < 1$, the *geometric series* $\sum_{n=0}^{\infty} \lambda^n$ is absolutely convergent, and its sum is

$$1 + \lambda + \lambda^2 + \lambda^3 + \dots + \lambda^n + \dots = \frac{1}{1 - \lambda}.$$

Dichotomy paradox. Using the above formula for the sum of the geometric series, you may try to convince [Zeno](#) that

$$1/2 + 1/4 + 1/8 + 1/16 + 1/32 + \dots = 1.$$

Decimal expansions. Also, you may convince yourself that $0.99999\dots$, which by definition is the sum of the series

$$\frac{9}{10} + \frac{9}{100} + \frac{9}{1000} + \frac{9}{10000} + \dots$$

is actually equal to 1. Moreover, you may learn how to recognize rational numbers as $0.3333\dots$ or $1.285714285714\dots$ from their periodic expansion. Indeed, a real number is rational if and only if its base 10 (or any other base $d \geq 2$) expansion is eventually periodic.

ex: Say if the following series are convergent, and, if so, compute their sum.

$$1 + 1/2 + 1/4 + 1/8 + 1/16 + \dots \quad 1 + 10 + 100 + 1000 + \dots \quad 1 + 1/10 + 1/100 + 1/1000 + \dots$$

$$\sum_{n=0}^{\infty} (4/5)^n \quad 9/10 + 9/100 + 9/1000 + \dots \quad 0.3333\dots$$

Convergence tests. Deciding convergence or divergence of a series is not easy. The only tool at our disposal is comparison with known series, and essentially the only known non-trivial series is the geometric one. Comparison means the obvious observation that $0 \leq x_n \leq y_n$ for any n sufficiently large implies the following two conclusions: $\sum_n y_n < \infty \Rightarrow \sum_n x_n < \infty$, and $\sum_n x_n = \infty \Rightarrow \sum_n y_n = \infty$.

Now, if $|x_n| \leq C \lambda^n$ for some constant $C > 0$ and any n sufficiently large, then the partial sums of the series $\sum_n x_n$ are bounded by a constant times the partial sums of the geometric series $\sum_n \lambda^n$, therefore the series $\sum_n x_n$ is absolutely convergent whenever $|\lambda| < 1$. This happens when $\limsup_{n \rightarrow \infty} |x_n|^{1/n} < 1$ (*root test*) or when $\limsup_{n \rightarrow \infty} |x_{n+1}/x_n| < 1$ (*ratio test*).

1.2 Linear recursions

Fibonacci model is the prototype of

Recursive linear equations. A *recursive linear equation* (or “finite difference linear equation”), a law

$$a_p x_{n+p} + a_{p-1} x_{n+p-1} + \cdots + a_1 x_{n+1} + a_0 x_n = f_n \quad (1.4)$$

which defines a sequence (x_n) given a set of “initial conditions” x_0, x_1, \dots, x_{p-1} and the known sequence (external force) f_n . Above, $a_0 \neq 0, a_1, \dots, a_{p-1}, a_p \neq 0$ are real or complex parameters. It is a discrete version of a linear ordinary differential equation of degree p with constant coefficients. We may interpret the left-hand side as Lx , obtained applying a linear operator L , obtained as a superposition of a finite number of powers D^k of the discrete derivative D , to the sequence $x = (x_n)$.

When $f_n = 0$ for all n , we get a *homogeneous recursive equation*

$$a_p x_{n+p} + a_{p-1} x_{n+p-1} + \cdots + a_1 x_{n+1} + a_0 x_n = 0. \quad (1.5)$$

or simply $Lx = 0$. The set of solutions of the homogeneous equation (1.5) is a vector space \mathcal{H} of dimension p , and the set of solutions of (1.4) is an affine space modeled on \mathcal{H} , i.e. has the form $(z_n) + \mathcal{H}$, where (z_n) is any (particular) solution of (1.4).

Eigenfunctions. The general recipe is: “linear homogeneous equations have exponential solutions”. The conjecture $x_n = z^n$ solves the recursive equation (1.5) if z is a root of the *characteristic polynomial*

$$P(z) = a_p z^p + a_{p-1} z^{p-1} + \cdots + a_1 z + a_0$$

In particular, if P has p distinct complex roots (which is the generic case), say z_1, z_2, \dots, z_p , then the general solution of the homogeneous equation is a linear combination

$$x_n = c_1 z_1^n + c_2 z_2^n + \cdots + c_p z_p^n$$

where the c_1, c_2, \dots, c_p are constants which depend on the initial conditions x_0, x_1, \dots, x_{p-1} .

ex: Find the general solution of the recurrence $x_{n+2} + 2x_{n+1} + x_n = 0$.

ex: Find an explicit formula for the Fibonacci numbers f_n ’s (which is known as *Binet’s formula*).

ex: Discuss what happens when the characteristic polynomial has non-simple roots (observe that if $z + \varepsilon$ and z are two roots, with $\varepsilon > 0$ small, then the superposition $((z + \varepsilon)^n - z^n) / \varepsilon \simeq n z^{n-1}$ is also a solution ...).

ex: Solve the discrete free particle Newton equation $D^2 x = 0$, i.e. $x_{n+2} - 2x_{n+1} + x_n = 0$.

ex: Consider the recursive equation

$$x_{n+2} = 2x_{n+1} + x_n.$$

Find the general solution. Find the solution with $x_0 = 0$ and $x_1 = 1$, and compute explicitly the first few terms of the sequence. Show that the quotients $q_n := x_{n+1}/x_n$ converge to $1 + \sqrt{2}$ when $n \rightarrow \infty$, and therefore

$$\frac{x_{n+1} - x_n}{x_n} \rightarrow \sqrt{2}$$

Obtain rational approximations of $\sqrt{2}$.

Generating functions. Given a sequence (x_n) , defined anyway, we may consider the (formal) power series

$$F(z) := \sum_{n \geq 0} x_n z^n$$

If the series has a non-zero radius of convergence (since the radius of convergence R is given by Hadamard formula $1/R = \limsup_{n \rightarrow \infty} \sqrt[n]{|x_n|}$, this happens when the x_n ’s grow at most exponentially, i.e. when $|x_n| \leq C\lambda^n$ for some $C > 0$ and $\lambda > 0$), it defines an analytic function $F(z)$

in some neighbourhood of the origin. Then, the original sequence may be recovered computing derivatives, since

$$x_n = \frac{F^{(n)}(0)}{n!}.$$

For this reason, $F(z)$ is called *generating function* of the sequence (x_n) .

You may find interesting the following characterization of rational functions.

Theorem 1.1. *A power series $\sum_{n \geq 0} x_n z^n$, converging in some neighbourhood of the origin, represents a rational function $F(z)$ iff the coefficients x_n satisfy a recursive linear homogeneous equation.*

Generating function of the Fibonacci numbers. If f_n denotes the n -th Fibonacci number, starting from $f_0 = f_1 = 1$, then the power series $\sum_{n \geq 0} f_n z^n$ represents the rational function

$$F(z) = \frac{1}{1 - z - z^2}$$

in a neighbourhood of the origin. Observe that it has a pole with smallest absolute value at $1/\phi$, and deduce that $\limsup_{n \rightarrow \infty} |f_n|^{1/n} = \phi$ (so that $f_n \sim \phi^n$, as we already knew).

ex: Give examples of sequences which do not satisfy any (finite) recursion.

Linear systems. A linear homogeneous recursive system is a law

$$x_{n+1} = Ax_n$$

for some vector valued sequence $x_n \in \mathbb{R}^k$, given a square matrix $A \in \text{Mat}_{k \times k}(\mathbb{R})$. The solution is

$$x_n = A^n x_0,$$

where $x_0 \in \mathbb{R}^k$ is the initial condition. The computation of powers A^n of a square matrix A is simplified if we can diagonalize it. For example, if the matrix has k distinct and real eigenvalues, then in the basis formed by the eigenvectors it is a diagonal matrix, say $A = \text{diag}(\lambda_1, \dots, \lambda_k)$, and its n -th power is simply the diagonal matrix $A^n = \text{diag}(\lambda_1^n, \dots, \lambda_k^n)$.

A finite difference equation of order p like

$$a_p y_{n+p} + a_{p-1} y_{n+p-1} + \dots + a_1 y_{n+1} + a_0 y_n = 0$$

is equivalent to a recursive linear homogeneous system $x_{n+1} = Ax_n$ for the vector values sequence $x_n := (y_n, y_{n-1}, \dots, y_{n-p+1})$.

ex: Write and solve the system which corresponds to Fibonacci problem.

1.3 Iteration of maps

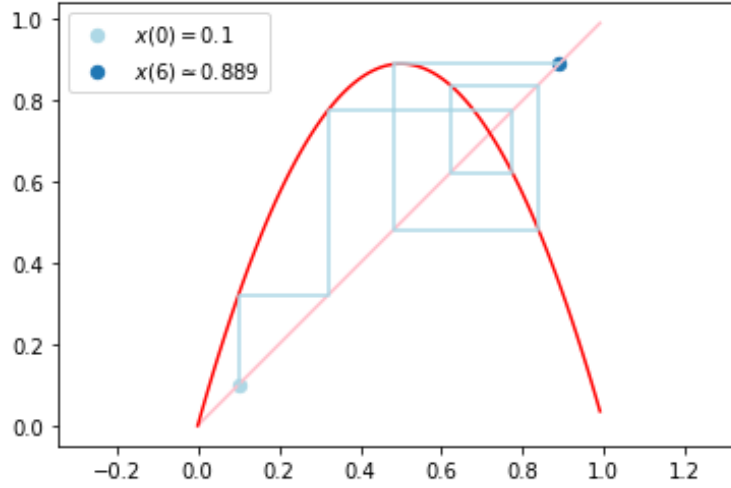
Iterations of maps. Given some space X (as the real line \mathbb{R} , an interval $I \subset \mathbb{R}$, an Euclidean space \mathbb{R}^N , and so on ...) and a transformation $f : X \rightarrow X$, we may form sequences according to

$$x_{n+1} = f(x_n)$$

given some initial condition x_0 . Such sequences are called *trajectories* of the map f .

Interval maps and cobweb plot. If X is an interval, we can follow trajectories using a “cobweb plot”: drawing vertical and horizontal lines connecting the points

$$(x, x) \mapsto (x, f(x)) \mapsto (f(x), f(x)) \mapsto (f(x), f^2(x)) \mapsto (f^2(x), f^2(x)) \mapsto (f^2(x), f^3(x)) \mapsto \dots$$



Cobweb plot of the quadratic map $f(x) = \lambda x(1 - x)$ when $\lambda = 3.56$.

Affine interval maps. As we have already seen, affine maps behave quite predictably. Indeed, the trajectories of an affine map like

$$f(x) = \lambda x + \alpha$$

with $\lambda \neq 1$, are sent, by the change of variable $y = x - \bar{x}$, where $\bar{x} = \alpha/(1 - \lambda)$ is the stationary solution, into the trajectories of $g(y) = \lambda y$, and the latter are geometric sequences. If $\lambda = 1$, trajectories are simply arithmetic series.

Nonlinearity. Non-linear recursive systems show much richer dynamics. Here is a short list of famous examples.

Hardy-Weinberg equilibrium. Consider the transmission of one gene with two alleles, say A and a . Let x_0, y_0, z_0 be the frequencies of the genotypes AA, Aa and aa , respectively, within some initial population. Then the probability to get the allele A or a in the formation of one gamete are

$$p_0 = x_0 + \frac{1}{2}y_0 \quad \text{e} \quad q_0 = 1 - p_0 = z_0 + \frac{1}{2}y_0,$$

respectively (so that $p_0 + q_0 = 1$). The offspring will therefore have genotypes AA, Aa or aa with frequencies

$$x_1 = p_0^2, \quad y_1 = 2p_0q_0 \quad \text{and} \quad z_1 = q_0^2.$$

(observe that $x_1 + y_1 + z_1 = p_0^2 + 2p_0q_0 + q_0^2 = (p_0 + q_0)^2 = 1$). The probabilities to get the allele A or a in the formation of one gamete in the second generation are

$$p_1 = x_1 + \frac{1}{2}y_1 \quad \text{and} \quad q_1 = z_1 + \frac{1}{2}y_1$$

Then an elementary computation (using only $p_0 + q_0 = 1$) shows that the second generation will have genotypes AA, Aa or aa with same frequencies as in the first generation, since

$$\begin{aligned} x_2 &= p_1^2 = \left(x_1 + \frac{1}{2}y_1\right)^2 = (p_0^2 + p_0q_0)^2 = p_0^2 = x_1 \\ y_2 &= 2p_1q_1 = 2\left(x_1 + \frac{1}{2}y_1\right)\left(z_1 + \frac{1}{2}y_1\right) = 2(p_0^2 + p_0q_0)(q_0^2 + p_0q_0) = 2p_0q_0 = y_1 \\ z_2 &= q_1^2 = \left(z_1 + \frac{1}{2}y_1\right)^2 = (q_0^2 + p_0q_0)^2 = q_0^2 = z_1 \end{aligned}$$

Thus, the distribution of the three genotypes attains a stationary value starting from the first generation (*Hardy⁴-Weinberg⁵ equilibrium/principle/law*).

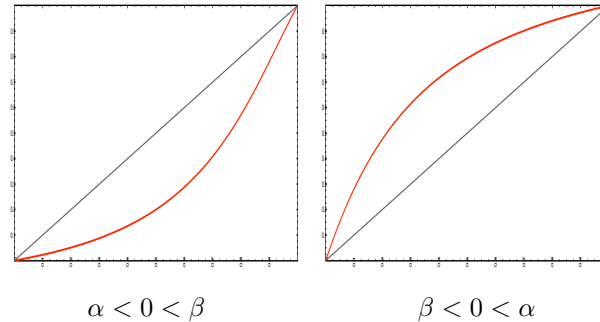
Fisher-Wright-Haldane model of natural selection. A simple model of natural selection was proposed by Fisher⁶, Wright⁷ and Haldane⁸ around 1930. It considers only one gene with two alleles, A and a . Biological success of the three different genotypes is modeled by certain “fitness coefficients” ϕ_{AA} , ϕ_{Aa} e ϕ_{aa} , which determine the different survival/reproduction rates. Let p_n and $q_n = 1 - p_n$ denote the frequencies of the alleles A and a , respectively, within the n -th generation. Then the frequency of the allele A at the next generation is

$$p_{n+1} = \frac{(1 + \alpha)p_n^2 + p_n q_n}{(1 + \alpha)p_n^2 + 2p_n q_n + (1 + \beta)q_n^2}$$

where we set $(1 + \alpha) = \phi_{AA}/\phi_{Aa}$ and $(1 + \beta) = \phi_{aa}/\phi_{Aa}$ (so that, since the fitness coefficients are positive by definition, both α and β are greater than -1).

The map $p_n \mapsto p_{n+1} = f(p_n)$ fully describes the time evolution of the population. It has two obvious fixed points, which are 0 and 1, and represent two homogeneous populations with only one allele.

If α and β have opposite signs (i.e. when the mixed genotype Aa has a fitness coefficient lying between the fitness coefficients of the pure genotypes AA and aa), these are the only fixed points. Observing at the graphs of $f(p)$ below



we see that if we start with any $0 < p_0 < 1$, then the sequence p_n converge to $p_n \rightarrow 0$ when $\alpha < 0 < \beta$ and converge to $p_n \rightarrow 1$ when $\beta < 0 < \alpha$. In both cases, the asymptotic population only contains the fittest allele, while the weakest get extincted.

More interesting things happen when α and β shares the same sign. The map $f(p)$ admits a third fixed point

$$\bar{p} = \frac{|\beta|}{|\alpha| + |\beta|},$$

strictly between 0 and 1, representing a mixed population.

When both α and β are positive (i.e. when both genotypes AA and aa perform better than Aa), then the equilibrium \bar{p} is unstable, a small perturbation $p_0 = \bar{p} \pm \varepsilon$ produces extinction of one of the two alleles, namely $p_n \rightarrow 0$ or 1, depending on the sign of the perturbation. This phenomenon is called *disruptive selection*.

When both α and β are negative (i.e. when the mixed genotype Aa is the fittest), then the equilibrium \bar{p} is stable, for any initial condition which is not 0 or 1 we get $p_n \rightarrow \bar{p}$. In particular, both alleles survive in the asymptotic population. This phenomenon is called *heterosis*.

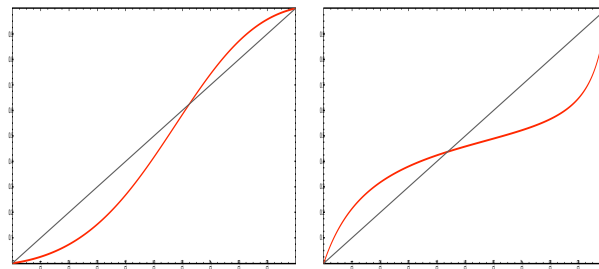
⁴G.H. Hardy, Mendelian proportions in a mixed population, *Science* **28** (1908), 49-50.

⁵W. Weinberg, Über den Nachweis der Vererbung beim Menschen, *Jahreshefte des Vereins für vaterländische Naturkunde in Württemberg* **64** (1908), 368-382.

⁶R.A. Fisher, *Genetical Theory of Natural Selection*, Clarendon 1930.

⁷S. Wright, Evolution in Mendelian populations, *Genetics* **16** (1931), 97-159.

⁸J.B.S. Haldane, A Mathematical Theory of Natural and Artificial Selection (1924-1934).



Disruptive selection: $0 < \alpha < \beta$. Heterosis: $\alpha < \beta < 0$.

The quadratic family. As soon as the interval map is not affine, trajectories are not easily understood. The simplest interval maps which are not affine are quadratic polynomials. A more realistic model of population dynamics in a limited environment seems to be

$$P_{n+1} = \lambda P_n (1 - P_n/M)$$

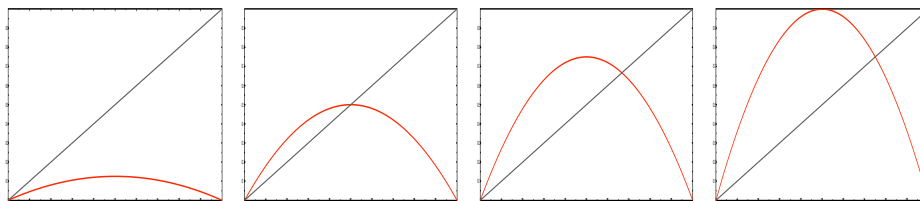
where the constant $M > 0$ is the maximal allowed population. Observe that $P_{n+1} < 0$ when $P_n > M$, which makes no physical sense (or may be we interpret it as “extinction”). The substitution $x_n = P_n/M$ transforms the above law into the adimensional law

$$x_{n+1} = \lambda x_n (1 - x_n),$$

The family of maps

$$f_\lambda(x) := \lambda x (1 - x) \tag{1.6}$$

is called *logistic map/transformation*⁹. The region where the relative population x_n makes (physical) sense is the unit interval $[0, 1]$ (which means real population between $0 \leq P_n \leq M$), and the map preserves the unit interval if the parameter ranges in the interval $0 \leq \lambda \leq 4$. Thus, we may think at f_λ as a map from the unit interval into itself. The particular map f_4 is also known as *von Neumann* or *Ulam map*.



Graphs of the logistic map when $\lambda = 0.5, 2, 3$ and 4 .

Stationary solutions are the trivial equilibrium 0 and the point $\bar{x} = (\lambda - 1)/\lambda$ (provided $\lambda > 1$). For small λ , the trajectories are previsible. As λ approaches 4 , they become quite wild.

When $\lambda > 4$, the unit interval is no longer preserved, and the map loses its physical/biological meaning. Nevertheless, it continues to be interesting for mathematicians.

ex: Write a code to simulate the system.

ex: Discuss what happens to trajectories when $0 < \lambda \leq 1$.

ex: Discuss what happens to trajectories when $1 < \lambda \leq 3$.

ex: Observe what happens when λ grows between 3 and 4 .

ex: What happens when $\lambda > 4$?

⁹Robert M. May, Simple mathematical models with very complicated dynamics, *Nature* **261** (1976), 459-467.

ex: Try to understand the dynamics of the following maps, defined in convenient intervals (some are easy, other are hard, if not impossible).

$$f(x) = \pm x^3 \quad f(x) = x^{1/3} \quad f(x) = x^3 \pm x$$

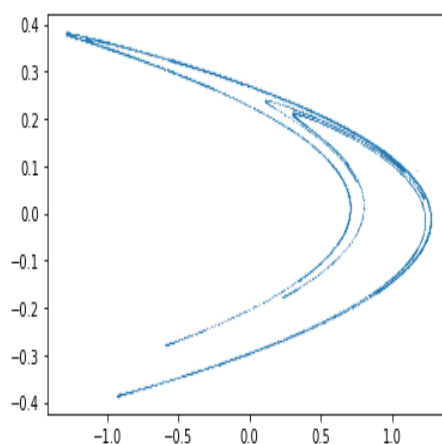
$$f(x) = x^2 + 1/4 \quad f(x) = |1 - x| \quad f(x) = x^2 - 2 \quad f(x) = \sin x \quad f(x) = \cos x$$

$$f(x) = x(1 - x) \quad f(x) = 2x(1 - x) \quad f(x) = 3x(1 - x) \quad f(x) = 4x(1 - x)$$

Henon map. The *Hénon map*¹⁰ is the map of the plane

$$\begin{cases} x_{n+1} = 1 + y_n - \alpha x_n^2 \\ y_{n+1} = \beta x_n \end{cases}$$

Depending on the values of its parameters, its trajectories show regular, “intermittent” or “chaotic” behavior. If you choose the parameters $\alpha \simeq 1.4$ and $\beta \simeq 0.3$, an initial condition like $x_0 \simeq 0.3$ and $y_0 \simeq 0.3$, and draw a sufficiently long orbit, you see the “Hénon attractor”



Hénon attractor.

1.4 Babylonians-Heron method to compute square roots

Searching for efficient methods to solve problems/equations is another source of interesting dynamical systems.

Babylonian-Heron algorithm. Consider the problem to find the side ℓ of a square given the value $a > 0$ of its area, i.e. to find the number which we call $\ell = \sqrt{a}$. A clever method, described by Heron¹¹, but probably already used by Babylonians^{12 13}, is as follows. We start with a rectangle with basis x_0 and height y_0 , “simple” numbers such that $x_0 y_0 = a$ (for example, if the area is an integer like $a = 2$, we may start with $x_0 = 3/2$ and $y_0 = 4/3$). We choose a second rectangle is such a way that its sides are nearer than the sides of the first rectangle. An obvious way to do it

¹⁰M. Hénon, A two-dimensional mapping with a strange attractor, *Comm. Math. Phys.* **50** (1976), 69-77.

¹¹“Since 720 has not its side rational, we can obtain its side within a very small difference as follows. Since the next succeeding square number is 729, which has 27 for its side, divide 720 by 27. This gives $26 \frac{2}{3}$. Add 27 to this, making $53 \frac{2}{3}$, and take half this or $26 \frac{5}{6}$. The side of 720 will therefore be very nearly $26 \frac{5}{6}$. In fact, if we multiply $26 \frac{5}{6}$ by itself, the product is $720 \frac{1}{36}$, so the difference in the square is $1/36$. If we desire to make the difference smaller still than $1/36$, we shall take $720 \frac{1}{36}$ instead of 729 (or rather we should take $26 \frac{5}{6}$ instead of 27), and by proceeding in the same way we shall find the resulting difference much less than $1/36$.”

Heron of Alexandria, *Metrica*, Book I.

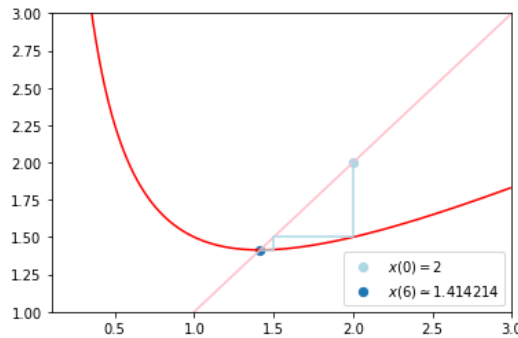
¹²C.B. Boyer, *A history of mathematics*, John Wiley & Sons, 1968.

¹³O. Neugebauer, *The exact sciences in antiquity*, Dover, 1969.

is to take as new basis the arithmetic mean $x_1 = (x_0 + y_0)/2$, which forces to take $y_1 = a/x_1$ as second height. And so on, if we are not satisfied yet. The recursion for that basis reads

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right).$$

Observe that if both a and the initial conjecture x_0 are rationals (the only numbers known to Babylonians), then all the x_n 's are also rationals.



First 6 iterations of the Heron method to find the square root of 2 starting from $x(0) = 2$.

Good rational approximations of $\sqrt{2}$. The algorithm converges, and quite fast. We could, as the Babylonians, put an initial guess $x_0 = 3/2$ for $\sqrt{2}$ (quite reasonable, since $1^2 < 2 < 2^2$), and find

$$x_1 = \frac{17}{12} \simeq 1.4166666666 \quad x_2 = \frac{577}{408} \simeq 1.41421568627 \quad x_3 = \frac{665857}{470832} \simeq 1.41421356237$$

As you see, the sequence stabilizes quite fast.

As a first attempt to explain this miracle, we could start looking at the recursive equations for the bases and the heights of the rectangles:

$$x_{n+1} = \frac{x_n + y_n}{2} \quad 1/y_{n+1} = \frac{1/x_n + 1/y_n}{2}$$

(so, the next height is the “harmonic mean” of the base and height). We see that the x_n 's and the y_n 's form decreasing and increasing sequences, respectively (disregarding the first guess, of course), namely

$$y_2 \leq y_3 \leq \dots \leq y_n \leq \dots \leq x_n \leq \dots \leq x_3 \leq x_2,$$

The real root is somewhere between, namely $y_n \leq \sqrt{a} \leq x_n$. Hence, we have an explicit control of the error. A computation shows that the lengths of those intervals, the differences $\varepsilon_n = x_n - y_n$ satisfy the recursion

$$\varepsilon_{n+1} < \frac{1}{2} \cdot \varepsilon_n$$

So, and initial “error” $\varepsilon_0 \leq 1$ (an easy achievement, since we easily recognize squares of integers) reduces to at least $\varepsilon_n \leq 2^{-n}$ after n iterations. The true error is actually much smaller. Indeed, in our example we may compute

$$\varepsilon_1 = \frac{17}{12} - \frac{24}{17} = \frac{1}{204} \simeq 0.005 \quad \text{and} \quad \varepsilon_2 = \frac{577}{408} - \frac{816}{577} = \frac{1}{235416} \simeq 0.000004$$

So that the first improved guess x_1 has already one correct decimals, and the second, x_2 has already four correct decimals!

What Babylonians didn't suspect is that if you start with a rational guess for $\sqrt{2}$, you get an infinite sequence of rational approximations, but the process never stops. This is due to

Theorem 1.2 (Pythagoras). *The square root of 2 is not rational.*

ex: A formula by Heron says that the area of a triangle with sides of lengths a , b and c , and semi-perimeter $s = (a + b + c)/2$, is given by

$$A = \sqrt{s(s-a)(s-b)(s-c)}$$

Estimate the area of a triangle with sides 7, 8 and 9.

ex: Estimate $\sqrt{13}$ with an error < 0.01 or 0.001 .

ex: Estimate how many iterations are necessary to obtain the first n correct decimals of $\sqrt{2}$ using Babylonians' method.

ex: Prove Pythagora's theorem 1.2 above (take a look at [HW59]).

Arithmetic-harmonic mean. Heron method can be better visualized as a bi-dimensional map. Given two positive numbers, x_0 and y_0 (the sides of a rectangle with area $a = x_0 y_0$), define recursively

$$(x_{n+1}, y_{n+1}) = f(x_n, y_n) := \left(\frac{x_n + y_n}{2}, \frac{2}{\frac{1}{x_n} + \frac{1}{y_n}} \right)$$

It is clear that the area function $A(x, y) := xy$ is preserved, i.e. $A(f(x, y)) = A(x, y)$. This means that trajectories belongs to hyperbolae $xy = \text{constant}$. Moreover, one easily sees that each trajectory $n \mapsto (x_n, y_n)$ converges to the diagonal, hence to the point (\sqrt{a}, \sqrt{a}) . For example, if we start with $(1, 2)$ we get $(\sqrt{2}, \sqrt{2})$ asymptotically.

Arithmetic-geometric mean. One is therefore tempted to generalize to other meaningful means. Given two positive numbers x and y , define recursively

$$a_{n+1} = \frac{1}{2}(a_n + g_n) \quad g_{n+1} = \sqrt{a_n g_n},$$

starting with $a_0 = (x + y)/2$ and $g_0 = \sqrt{xy}$, the arithmetic and the geometric mean of x and y , respectively. The arithmetic-geometric mean inequality (the fact that $(x + y)^2 \geq 0$) says that $g_n \leq a_n$, and therefore

$$g_{n+1} = \sqrt{a_n g_n} \geq \sqrt{g_n g_n} = g_n$$

Since both sequences a_n and g_n are between the minimum and the maximum of x and y , this implies that g_n converges, to some (positive) limit p . The sequence a_n also converges, and to the same limit, since

$$a_n = g_{n+1}^2 / g_n \rightarrow p$$

The common limit is called *arithmetic-geometric mean* of x and y , say $p =: \text{AGM}(x, y)$. What is not trivial is a formula for the limit, and this is due to Gauss: it says that

$$\text{AGM}(x, y) = \frac{\pi}{4} \frac{x + y}{K\left(\frac{x-y}{x+y}\right)}$$

where

$$K(k) := \int_0^{\pi/2} \frac{d\theta}{\sqrt{1 - k^2 \sin^2 \theta}}$$

is the “complete elliptic integral of the first kind”.

Side and diagonal numbers. Define *side* and *diagonal* numbers recursively as

$$s_{n+1} = s_n + d_n \quad d_{n+1} = 2s_n + d_n$$

respectively, starting from $s_0 = d_0 = 1$. A computation shows that

$$d_n^2 - 2s_n^2 = \pm 1$$

where the sign depends on the parity of n . Thus, when n is odd, we get a whole family of integer solutions of the Pell equation $x^2 - 2y^2 = 1$. Since both s_n and d_n grow, the ratio d_n/s_n tends to $\sqrt{2}$ when $n \rightarrow \infty$, with an error of size $1/\sqrt{s_n}$.

1.5 From Newton method to Julia and Fatou sets

Finding \sqrt{a} means solving the polynomial equation $z^2 - a = 0$. What about finding roots of a generic polynomial ?

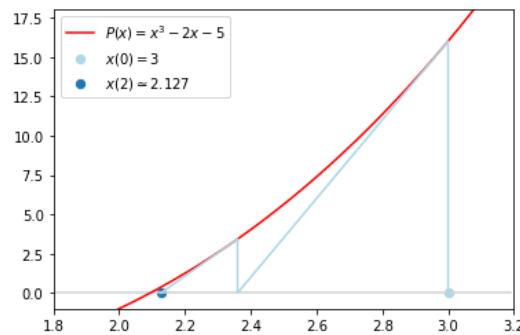
Newton-Raphson iterative scheme. “Newton method” is a method proposed by Joseph Raphson around 1690 to approximate roots of a polynomial $p(x)$ (Newton used it to solve $x^3 - 2x - 5 = 0$). It consists in starting with an initial conjecture x_0 near to some root, and then improve it using the linear approximation

$$p(x) \simeq p(x_0) + p'(x_0)(x - x_0).$$

This idea leads to the recursion

$$x_{n+1} = x_n - \frac{p(x_n)}{p'(x_n)}.$$

It is clear that if the sequence converges, i.e. $x_n \rightarrow x_\infty$, and if $p'(x_\infty) \neq 0$, then the limit x_∞ is a root.



Search for a root of $x^3 - 2x - 5$ using Newton iterations.

ex: Use Newton method to solve Newton’s problem, i.e. find the roots of $x^3 - 2x - 5$.

ex: Show that Newton method to solve $x^2 - a = 0$ corresponds to babylonian-Heron iterative scheme.

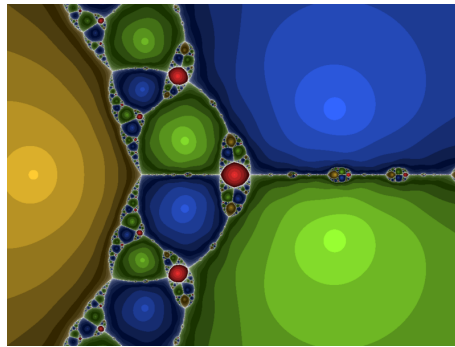
ex: Use Newton method to approximate the Greeks’ *ratio*, the positive root of $x^2 - x - 1$. Then, compare with the babylonian-Heron method (i.e., estimate $\sqrt{5}$, then sum 1 and divide by 2).

ex: Write and implement Newton method to find n -th roots, i.e. to solve $x^n - a = 0$.

Newton’s fractals. In 1879 Cayley observed that the above method could be also used to approximate complex roots of complex polynomials $p(z) \in \mathbb{C}[z]$. It amounts to iterate the rational function

$$f(z) = z - \frac{p(z)}{p'(z)}$$

The problem is therefore to understand when, i.e for which initial values z_0 , the sequence z_n converges to one of the roots. The “basins of attraction” of the different roots draw beautiful and unexpected patterns in the complex plane.



Basins of attraction of the roots of $2z^3 - 2z + 2$ in \mathbb{C}
(from http://en.wikipedia.org/wiki/Newton_fractal).

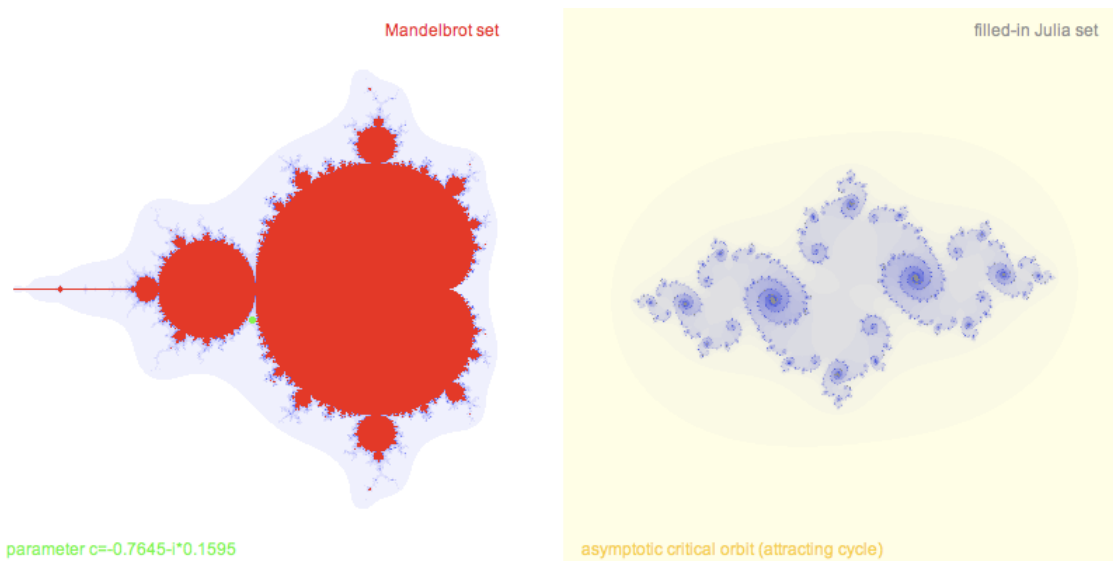
Iteration of rational functions in the Riemann sphere. The natural generalization is to take a rational function of a complex variable $f(z) \in \mathbb{C}(z)$, which is an endomorphism of the Riemann sphere $\overline{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$, i.e. try to understand its trajectories, i.e. the iteration $z_{n+1} = f(z_n)$.

Most studied is iteration of the family of quadratic polynomials

$$f(z) = z^2 + c$$

depending on a parameter $c \in \mathbb{C}$. Its beauty was foreseen by Gaston Julia¹⁴ and Pierre Fatou¹⁵ at the beginning of the XX century, revealed with the help of the first modern computers by Benoît Mandelbrot, and then studied by a variety of great mathematicians (like Adrian Douady, Dennis Sullivan, John Milnor, Misha Lyubich, Jean-Christophe Yoccoz, Curtis McMullen, ...) starting from the 80's of the last century.

Pictures of the Mandelbrot and Julia sets. Below, you may find a picture of what Julia and Fatou could only dream about.



Mandelbrot set (left) and Julia set (right) of the polynomial $z^2 + c$ with $c \simeq -0.7645 - i \cdot 0.1595$.

(from <http://w3.math.uminho.pt/~scosentino/bestiario/julia.html>)

The red hearts on the left form the *Mandelbrot set*, the set of those values of the parameter c such the orbit of critical point $z_0 = 0$ is bounded. The almost invisible grey points on the right

¹⁴G. Julia, Mémoire sur l'iteration des fonctions rationnelles, *Journal de Mathématiques Pures et Appliquées*, **8** (1918), 47-245.

¹⁵P. Fatou, Sur les substitutions rationnelles, *Comptes Rendus de l'Académie des Sciences de Paris*, **164** (1917) 806-808, and **165** (1917), 992-995.

form the *filled-in Julia set*, the set of initial values z_0 with bounded orbits (once fixed a value of c). Blue colors, which help to see the Julia set, are chosen depending on the speed with which other trajectories diverge to ∞ .

Much more beautiful pictures, and then movies and so on, may be found in this page by Jos Leys: <http://www.josleys.com>

2 Differential equations and flows

2.1 Flows

The main way in which dynamical systems enter in physics is through differential equations.

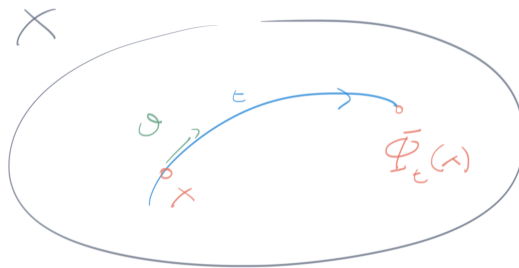
Flows of vector fields. Let X be a differentiable manifold (as, for example, an open region of \mathbb{R}^N), and let v be a vector field on X . If we assume that the autonomous differential equation

$$\dot{x} = v(x)$$

with any given initial condition $x(0) = x$, has solutions $t \mapsto x(t)$ which exist for any time $t \in \mathbb{R}$ (as is the case when v is smooth and X is compact), then the *flow* of the vector field v is the action $\Phi : \mathbb{R} \times X \rightarrow X$ given by $\Phi_t(x) = x(t)$. Indeed, it is clear that Φ_0 is the identity map, and that

$$\Phi_t \circ \Phi_s = \Phi_{t+s}$$

for any $t, s \in \mathbb{R}$. Therefore, $\Phi_{-t} = (\Phi_t)^{-1}$.



Conversely, given a one-parameter group of diffeomorphisms Φ_t , one defines the phase velocity according to

$$v(x) := \left. \frac{d}{dt} \Phi_t(x) \right|_{t=0} = \lim_{t \rightarrow 0} \frac{\Phi_t(x) - x}{t}$$

The group property then implies that the curve $t \mapsto x(t) = \Phi_t(x)$ satisfies

$$\dot{x}(t) = \lim_{s \rightarrow 0} \frac{\Phi_{t+s}(x) - \Phi_t(x)}{s} = \lim_{s \rightarrow 0} \frac{\Phi_s(\Phi_t(x)) - \Phi_t(x)}{s} = v(x(t))$$

and therefore is a solution of the autonomous differential equation $\dot{x} = v(x)$ with initial condition $x(0) = x$.

Also interesting are *semi-flows* Φ_t , which are defined only for non-negative times $t \geq 0$.

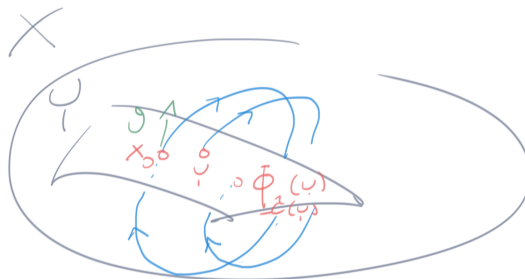
A flow or semi-flow is called *continuous time dynamical system*, and indeed our basic definitions in the previous chapter are adaptations of physicists' ideas and terminology about flows of vector fields. The map $t \mapsto \Phi_t(x)$ is called *trajectory* of the (initial) point x , and its image $O^+(x) = \{\Phi_t(x) : t \in \mathbb{R}_+\}$ is called (*forward*) *orbit* of x . If it happens, as usual in classical mechanics, that flows are defined for all times $t \in \mathbb{R}$, then the set $O(x) = \{\Phi_t(x) : t \in \mathbb{R}\}$ is called *orbit* of x .

From flows to maps, discretization. Given a flow Φ_t on X , one could specialize to discrete time looking at the system at multiples integers $n\tau$ of a given time-unit $\tau > 0$, and this amounts to iterate the transformation $f = \Phi_\tau$.

More interesting is the following construction.

Poincaré maps. Let Φ_t be the flow of the autonomous differential equation $\dot{x} = v(x)$ on a manifold X , and let $Y \subset X$ be a submanifold of codimension one which is transversal to the flow (i.e. the tangent space $T_x Y$ does not contain the vectors $v(x)$ for any $x \in Y$).

If $x_0 \in Y$ is a periodic point, say $\Phi_\tau(x_0) = x_0$ for some period $\tau > 0$, then nearby points $x \in Y$ also return to Y after some time near to τ . Thus, one could define, in a sufficiently small neighbourhood $U \subset Y$ of x_0 , a *first return/Poincaré map* $f : U \rightarrow Y$, sending a point $x \in U$ into $\Phi_{\tau(x)}(x)$ if $\tau(x)$ is the smallest positive time $t > 0$ such that $\Phi_t(x) \in Y$. This construction is even possible around a point which is not periodic, provided its orbit returns to Y sufficiently near.



Moreover, it may be also happens that the flow allows a *global (Poincaré) section*, a codimension one submanifold $Y \subset X$ transversal to the vector field v such that the orbit of any point $y \in Y$ eventually returns to Y after a minimal time

$$\tau(y) := \inf\{t > 0 \text{ s.t. } \Phi_t(y) \in Y\} < \infty,$$

called *first return time*. This allows to consider a globally defined *first return/Poincaré map* $f : Y \rightarrow Y$, according to

$$f(y) := \Phi_{\tau(y)}(y).$$

Linear flows on the two-torus and rotations of the circle. A constant vector field $v = (a, b)$ generates a linear flow

$$\Phi_t : (x, y) \mapsto (x + at, y + bt)$$

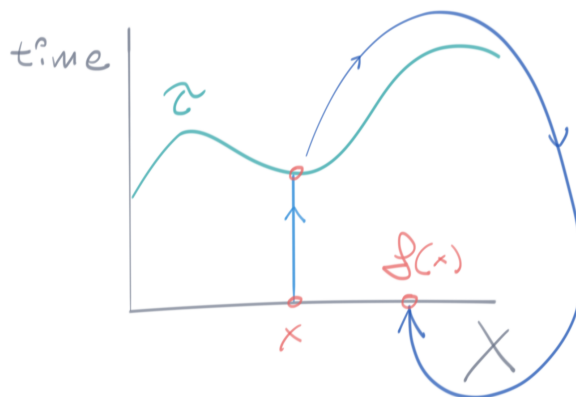
on the plane. This flow is clearly invariant under translations by integer vectors, and therefore it defines a flow Φ_t on the two dimensional torus $\mathbb{T}^2 := \mathbb{R}^2/\mathbb{Z}^2$. The circle $\mathbb{R}/\mathbb{Z} \simeq C := \{(x + \mathbb{Z}, 0 + \mathbb{Z})\} \subset \mathbb{T}^2$ is transversal to the vector field if $b \neq 0$. The orbit of a point $(x, 0) \in C$ goes back to the section after a time $\tau = 1/b$ (if $b > 0$) to the point $\Phi_\tau(x, y) = (x + a/b, 0)$. Thus, the first return map is $f : x + \mathbb{Z} \mapsto x + \alpha + \mathbb{Z}$, a rotation of the circle by an “angle” $\alpha = a/b$.

Suspension flows. Poincaré construction of a first return map out of a flow admits an inverse. Given a map $f : X \rightarrow X$, one can define the *mapping torus* X_f as the Cartesian product $X \times [0, 1]$, with coordinates (x, t) with $x \in X$ and $t \in [0, 1]$, modulo the equivalence relation $(x, 1) \sim (f(x), 0)$. The flow of the vertical vector field $\partial/\partial t$ on X_f (which is a smooth manifold if X is) is called *suspension* of f . It is clear that it admits a global Poincaré section $X \times \{0\} \simeq X$, and its first return map is precisely f .

More generally, given a map $f : X \rightarrow X$ and a *roof function* $\tau : X \rightarrow \mathbb{R}_+$ bounded away from 0, one can consider the space $X_{f,\tau} = Y/\sim$ obtained as

$$Y = \{(x, t) : x \in X, 0 \leq t \leq \tau(x)\}$$

modulo the equivalence relation $(x, \tau(x)) \sim (f(x), 0)$. The flow of the vertical vector field $\partial/\partial t$ on X_f is called *suspension* of f with *height* τ . Again, it admits a global Poincaré section $X \times \{0\} \simeq X$, and its first return map is f .



2.2 Structure of physical models

Classical mechanics is the natural source of interesting dynamical systems.

Newtonian mechanics. According to greeks, the “velocity” $\dot{q} = \frac{d}{dt}q$ of a planet, where $q \in \mathbb{R}^3$ is its position in our Euclidean space and t is time, was determined by gods or whatever forced planets to move around circles. Then came Galileo, and showed that gods could at most determine the “acceleration” $\ddot{q} = \frac{d^2}{dt^2}q$, since the laws of physics should be written in the same way by an observer in any reference system at uniform rectilinear motion with respect to the fixed stars. Finally came Newton, who decided that what gods determined was to be called “force”, and discovered that the trajectories of planets, fulfilling Kepler’s experimental three laws¹⁶, were solutions of his famous (second order differential) equation

$$m\ddot{q} = F$$

where m is the mass of the planet, and where the attractive force F between the planet and the Sun is proportional to the product of their masses and inverse proportional to the square of their distance.

Later, somebody noticed that most observed forces were “conservative”, could be written as $F = -\nabla V$, for some real valued function $V(q)$ called “potential energy”. There follows that Newton equations can be written as $m\ddot{q} = -\nabla V$, and that the “total energy”

$$E = \frac{1}{2}m\|\dot{q}\|^2 + V(q)$$

is constant along trajectories. The function $\frac{1}{2}m\|\dot{q}\|^2$ is called “kinetic energy” of the system.

Lagrangian and variational principle. An alternative (and indeed useful) formulation of Newtonian mechanics is the one developed by Lagrange. He defined (what we now call) the “Lagrangian” of the system as the difference between the kinetic energy and the potential energy

$$L(q, \dot{q}) = \frac{1}{2}m\|\dot{q}\|^2 - V(q)$$

and observed that Newton equations are equivalent to the (Euler)-Lagrange equations

$$\frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}} \right) = \frac{\partial L}{\partial q}$$

This is important because solutions of the Euler-Lagrange equations are critical points of the *action*

$$S[q(t)] = \int_{t_0}^{t_1} L(q(t), \dot{q}(t)) dt.$$

Thus, we may look at trajectories of a physical system as (local) minimizers of a certain variational problem. This often allows to find the trajectories without even solving the equations of motion.

Hamiltonian mechanics. The product $p = m\dot{q} = \partial L / \partial \dot{q}$ is called “(linear) momentum”, and, since p/m is the gradient of the kinetic energy $K(p) = \|p\|^2/2m$, Hamilton could write Newton’s second order differential equations as the system of first order differential equations

$$\dot{q} = \frac{\partial H}{\partial p} \quad \dot{p} = -\frac{\partial H}{\partial q}$$

¹⁶In *Astronomia nova*, 1609, and *Harmonices mundi*, 1619, Johannes Kepler published his three laws of planetary motions:

- i) planets moves in ellipses with focus at the Sun,
- ii) the radius vector describes equal areas in equal times,
- iii) the squares of the periods are to each other as the cubes of the mean distance from the Sun.

It was with the purpose to derive Kepler laws from a second order differential equation $m\ddot{q} = F$ that Isaac Newton realized that the force of gravitational attraction between the Sun and a planet (hence between any two bodies!) should be proportional to m/ρ^2 (*Philosophiae naturalis principia mathematica*, 1687).

where $H(q, p) = K(p) + V(q)$ is the total energy as function of q and p , nowadays called “Hamiltonian”. It is a simple check that the energy is a constant of the motion, since

$$\frac{d}{dt}H = \frac{\partial H}{\partial q} \cdot \dot{q} + \frac{\partial H}{\partial p} \cdot \dot{p} = \frac{\partial H}{\partial q} \cdot \frac{\partial H}{\partial p} - \frac{\partial H}{\partial p} \cdot \frac{\partial H}{\partial q} = 0$$

The modern abstract formulation of classical mechanics is as follows. Let (X, ω) be a symplectic manifold, i.e. a differentiable manifold X of even dimension $2n$, equipped with a smooth closed differential two-form ω such that $\omega^n \neq 0$. Darboux theorem says that locally one can choose “canonical” coordinates $(q_1, \dots, q_n, p_1, \dots, p_n)$ such that $\omega = \sum_{k=1}^n dp_k \wedge dq_k$. The standard example is the cotangent bundle $T^*\mathbb{R}^N$ of the Euclidean vector space \mathbb{R}^N , whose coordinates are positions q_k and momenta p_k .

Let $H : X \rightarrow \mathbb{R}$ be a smooth function, called “Hamiltonian” and thought as the “energy” of the system. Typically, it has the form “kinetic energy+potential energy”, where the kinetic energy is a positive definite quadratic form in the momenta p , and the potential energy is a function V depending on the positions q and possibly on the momenta p . The Hamiltonian vector field v is defined by the identity $dH = i_v \omega$, and the *Hamiltonian flow* is the flow of v . In canonical coordinates, the equations of motion read

$$\dot{q}_k = \frac{\partial H}{\partial p_k} \quad \dot{p}_k = -\frac{\partial H}{\partial q_k}$$

It happens that the Hamiltonian flow Φ_t preserves the energy, namely $H(\Phi_t(x)) = H(x)$ for any $x \in X$ and any time $t \in \mathbb{R}$, as follows from the fact that $\mathcal{L}_v H = 0$.

Also, according to Liouville theorem, it preserves the volume form ω^n , defined in canonical coordinates by the volume element $dq_1 \wedge \dots \wedge dq_n \wedge dp_1 \wedge \dots \wedge dp_n$. In particular, if the phase space is compact, it preserves a probability measure.

Free motion. Free motion in an inertial frame is described by the Lagrangian $L = \frac{1}{2}m\|\dot{q}\|^2$. The equations of motion are

$$\ddot{q} = 0.$$

Solutions are straight lines $q(t) = c + vt$, for same initial position $q(0) = c$ and velocity $\dot{q}(0) = v$.

Free fall. Free fall near the Earth’s surface is modeled by the Lagrangian $L = \frac{1}{2}m\|\dot{q}\|^2 - mgz$, where $g \simeq 9.8 \text{ m/s}^2$ is the gravitational acceleration and z is the height of the particle (assumed much smaller than the Earth’s diameter), the third coordinates of $q = (x, y, z)$. The equation of motion for the height is

$$\ddot{z} = g.$$

Solution are parabola $z(t) = c + vt - \frac{1}{2}gt^2$, for some initial height $c = z(0)$ and some initial velocity $v = \dot{z}(0)$.

Geodesic flows. The simplest mechanical system, the free motion of a particle, belongs to the class of geodesic flows. Let (M, g) be a Riemannian manifold, g being the Riemannian metric. Let SM be the unit tangent bundle of M . If M is geodesically complete, to every unit vector $v \in SM$ there corresponds a unique geodesic line (i.e. a local isometry) $c : \mathbb{R} \rightarrow M$ such that $\dot{c}(0) = v$. The *geodesic flow* is the action $\Phi : \mathbb{R} \times SM \rightarrow SM$, defined as $\Phi_t(v) = \dot{c}(t)$.

Particularly interesting are geodesic flows over homogeneous spaces. Apart from the rather trivial example of flat spaces, a source of interesting dynamical properties is the geodesic flow on a manifold with constant negative curvature. The prototype is as follows. The group $G = PSL(2, \mathbb{R})$ can be seen as the orientation preserving isometry group of the Poincaré half-plane \mathbb{H} , equipped with the hyperbolic metric of sectional curvature -1 . Its action is transitive. Since the stabilizer of a point in the half-plane is isomorphic to the group of rotations $SO(2)$, we can identify $S\mathbb{D}$ with G . Now, let Γ be a discrete cocompact subgroup of G with no torsion. The quotient space $\Sigma = \mathbb{D}/\Gamma$ is a compact Riemann surface, which comes equipped with a Riemannian metric of sectional curvature -1 , and its unit tangent bundle is diffeomorphic to G/Γ . The geodesic flow on $S\Sigma$ is then the algebraic flow $\Phi : \mathbb{R} \times G/\Gamma \rightarrow G/\Gamma$ defined as $\Phi_t(g\Gamma) = e_t g\Gamma$, where

$$e_t = \begin{pmatrix} e^{t/2} & 0 \\ 0 & e^{-t/2} \end{pmatrix}$$

2.3 Integration of one-dimensional systems

Some techniques to integrate ordinary differential equations (ODEs) like $\dot{x} = v(x, t)$ when the phase space is one or two-dimensional.

Integrating simple ODEs. The simplest case occurs when the velocity field v does not depend on the phase space variable x , hence

$$\dot{x} = v(t),$$

where $v(t)$ is some given (piecewise) continuous function of time. This just says that x must be a primitive of v , and the fundamental theorem of calculus (i.e. Leibniz and/or Newton's discovery) tells us how to compute such a primitive:

$$x(t) = x_0 + \int_{t_0}^t v(s) ds.$$

Here you may observe that this class of ODEs have “symmetries”. The line field does not depend on x , hence slopes of solutions are the same along horizontal lines ($t = \text{constant}$) in the extended phase space $X \times \mathbb{R}$. There follows that any translate $\varphi(t) + c$ of a solution $\varphi(t)$ is still a solution.

Autonomous first order ODEs and their flows. A first order ODE of the form

$$\dot{x} = v(x),$$

where the velocity field v does not depend on time, is called *autonomous*. Most fundamental equations of physics (those describing closed systems, without external forces) can be written as autonomous first order ODEs, and this corresponds to time-invariance of physical laws.

Here you may notice symmetries again. The line field v of an autonomous equation is constant along vertical lines ($x = \text{constant}$) of the extended phase space $X \times \mathbb{R}$. Hence any translate $\varphi(t+s)$ of a solution $\varphi(t)$ is still a solution. This is the manifestation of time-invariance of a law codified by an autonomous ODE. This also implies that there is no loss of generality in restricting to an initial time $t_0 = 0$.

Equilibrium solutions. First, we observe that an autonomous equation may admit constant solutions. Indeed, if x_0 is a *singular point* of the vector field v , i.e. a point where $v(x_0) = 0$, then the constant function

$$x(t) = x_0 \quad \forall \quad t \in \mathbb{R}$$

obviously solves the equation. Such solutions, which do not change with time, are called *equilibrium*, or *stationary*, solutions.

Solutions near non-singular points. The trick used to “guess” other solutions, when the phase space is one-dimensional, i.e. $X \subset \mathbb{R}$, is a first instance of the method of “separation of variables”. Fix a *non-singular point* of the velocity field, i.e. a point x_0 where $v(x_0) \neq 0$. We want to solve the Cauchy problem with initial condition $x(t_0) = x_0$. First, rewrite the equation $dx/dt = v(x)$ formally as “ $dx/v(x) = dt$ ” (multiply by dt and divide by $v(x)$, so that all x 's are on the left and all t 's are on the right). Instead of trying to make sense to this last expression (which is possible, of course, and here you can appreciate the beauty of Leibniz' notation dx/dt for derivatives!), observe that it is suggesting that $\int dx/v(x) = \int dt$. Now assume that the velocity field v is continuous and let $J = (x_-, x_+)$ be the maximal interval containing x_0 where v is different from zero. Integrating, from x_0 to $x \in J$ on the left and from t_0 to t on the right, we obtain a differentiable function $x \mapsto t(x)$ defined as

$$t(x) = t_0 + \int_{x_0}^x \frac{dy}{v(y)}$$

for any $x \in J$. Now, observe that the derivative dt/dx is equal to $1/v$. Since, by continuity, $1/v$ does not change its sign in J , our $t(x)$ is a strictly monotone continuously differentiable function. We can invoke the inverse function theorem and conclude that the function $t(x)$ is invertible. This

prove that the above relation defines actually a continuously differentiable function $t \mapsto x(t)$ in some interval $I = t(J)$ of times around t_0 . Finally, you may want to check that the function $t \mapsto x(t)$ solves the Cauchy problem: just compute the derivative (using the inverse function theorem),

$$\begin{aligned}\dot{x}(t) &= 1 / \left(\frac{dt}{dx}(x(t)) \right) \\ &= v(x),\end{aligned}$$

and check the initial condition. Observe that the function $t(x) - t_0$ has then the interpretation of the “time needed to go from x_0 to x ”.

At the end of the story, if you are lucky enough and know how to invert the function $t(x)$, you’ll get an explicit solution as

$$x(t) = F^{-1}(t - t_0 + F(x_0)),$$

where F is any primitive of $1/v$. Close inspection of the above reasoning shows that the local solution you’ve found is indeed the unique one. Namely, we have the following

Theorem 2.1. *Let $v(x)$ be a continuous velocity field and let x_0 be a non-singular point of v . Then there exist one and only one solution of the Cauchy problem $\dot{x} = v(x)$ with initial condition $x(t_0) = x_0$ in some sufficiently small interval I around t_0 . Moreover, the solution $x(t)$ is the inverse function of*

$$t(x) = t_0 + \int_{x_0}^x \frac{dy}{v(y)},$$

defined in some small interval J around x_0 .

Proof. Here we give the pedantic proof. Let J be as above. Define a function $H : \mathbb{R} \times J \rightarrow \mathbb{R}$ as

$$H(t, x) = t - t_0 - \int_{x_0}^x \frac{dy}{v(y)}.$$

If $t \mapsto \varphi(t)$ is a solution of the Cauchy problem, then computation shows that $\frac{d}{dt}H(t, \varphi(t)) = 0$ for any time t . There follows that H is constant along the solutions of the Cauchy problem. Since $H(t_0, x_0) = 0$, we conclude that the graph of any solution belongs to the level set $\Sigma = \{(t, x) \in \mathbb{R} \times J \text{ s.t. } H(t, x) = 0\}$. Now observe that H is continuously differentiable and that its differential $dH = dt + dx/v(x)$ is never zero. Actually, both partial derivatives $\partial H/\partial t$ and $\partial H/\partial x$ are always different from zero. Hence we can apply the implicit function theorem and conclude that the level set Σ is, in some neighbourhood $I \times J$ of (t_0, x_0) , the graph of a unique differentiable function $x \mapsto t(x)$, as well as the graph of a unique differentiable function $t \mapsto x(t)$, the inverse of t , which as we have already seen solves the Cauchy problem. \square

On the failure of uniqueness near singular points. The interval $I = t(J)$ where the solution is defined need not be the entire real line: solutions may reach the boundary of J , i.e. one of the singular points x_{\pm} of the velocity field, in finite time. Since singular points are themselves equilibrium solutions, this implies that solutions of the Cauchy problem at singular points may not be unique, under such mild conditions (continuity) for the velocity field. Later we’ll see Picard’s theorem, which prescribes stronger regularity conditions on the velocity field v under which the Cauchy problem admits unique solutions for any initial condition in the extended phase space.

Counter-example. Both curves $x(t) = 0$ and $x(t) = t^3$ solve the equation

$$\dot{x} = 3x^{2/3}$$

with initial condition $x(0) = 0$. The problem here is that the velocity field $v(x) = 3x^{2/3}$, although continuous, is not differentiable and not even Lipschitz at the origin. You may notice that the solution starting, for example, at $x_0 = 1$ reaches (or better comes from) the singular point $x_- = 0$ in finite time, since

$$t(x_-) - t(x_0) = \int_1^0 \frac{1}{3} y^{-2/3} dy = -1.$$

One-dimensional Newtonian motion in a time independent force field. The one-dimensional motion of a particle of mass m subject to a force $F(x)$ that does not depend on time is described by the Newton equation

$$m\ddot{x} = -U'(x),$$

where the potential $U(x) = -\int F(x)dx$ is some primitive of the force. The total energy

$$E(x, \dot{x}) = \frac{1}{2}m\dot{x}^2 + U(x)$$

(which of course is defined up to an arbitrary additive constant) of the system is a constant of the motion, i.e. is constant along solutions of the Newton equation. In particular, once a value E of the energy is given (depending on the initial conditions), the motion takes place in the region where $U(x) \leq E$, since the kinetic energy $\frac{1}{2}m\dot{x}^2$ is non-negative. Conservation of energy allows to reduce the problem to the first order ODE

$$\dot{x}^2 = \frac{2}{m}(E - U(x)),$$

which has the unpleasant feature to be quadratic in the velocity \dot{x} . Meanwhile, if we are interested in a one-way trajectory going from some x_0 to x , say with $x > x_0$, we may solve for \dot{x} and find the first order autonomous ODE

$$\dot{x} = \sqrt{\frac{2}{m}(E - U(x))}.$$

There follows that the time needed to go from x_0 to x is

$$t(x) = \int_{x_0}^x \frac{dy}{\sqrt{\frac{2}{m}(E - U(y))}}.$$

The inverse function of the above $t(x)$ will give the trajectory $x(t)$ with initial position $x(0) = x_0$ and initial positive velocity $\dot{x}(0) = \sqrt{\frac{2}{m}(E - U(x_0))}$, at least for sufficiently small times t .

The exponential. The exponential function, according to Walter Rudin “the most important function in mathematics” ([Ru87], 1st line of page 1), is the unique solution of the autonomous differential equation

$$\dot{x} = x$$

with initial condition $x(0) = 1$. If we try a power series like $a_0 + a_1t + a_2t^2 + a_3t^3 + \dots$, the differential equation gives the recursion $na_n = a_{n-1}$ for the coefficients, while the initial condition yields $a_0 = 1$. Thus, the solution is $x(t) = 1 + t + t^2/2 + t^3/6 + \dots$.

Actually, it is convenient to complexify time, i.e. take $z = t + i\theta \in \mathbb{C}$ with $t, \theta \in \mathbb{R}$, and define the *exponential* as the power series

$$\exp(z) := 1 + z + \frac{z^2}{2} + \frac{z^3}{6} + \frac{z^4}{24} + \dots = \sum_{n \geq 0} \frac{z^n}{n!}$$

Since $\limsup_{n \rightarrow \infty} (1/n!)^{1/n} = 0$, the radius of convergence is ∞ , hence the power series defines an entire function, i.e. a holomorphic function $\exp : \mathbb{C} \rightarrow \mathbb{C}$. Deriving each term of the series, we easily verify that indeed $\exp' = \exp$. The initial condition $\exp(0) = 1$ is obvious. From absolute convergence of the series and algebraic manipulation we also get the group property

$$\exp(z + w) = \exp(z) \cdot \exp(w)$$

for any $z, w \in \mathbb{C}$, saying that \exp is a homomorphism of the additive group \mathbb{C} into the multiplicative group $\mathbb{C}^\times = \mathbb{C} \setminus \{0\}$. In particular, $\exp(-z) = 1/\exp(z)$, so that the exponential $\exp(z)$ is never 0. This also justifies our notation $\exp(z) = e^z$, where

$$e := \exp(1) = 1 + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots \simeq 2.7182818284590452353602874713526624977572\dots$$

(another famous irrational, actually a transcendental number!). For real time $z = t$, we recover the familiar model of “exponential growth” $t \mapsto e^t$, a strictly increasing function from the additive group \mathbb{R} onto the multiplicative group $\mathbb{R}_+ =]0, \infty[$, growing faster than any power t^n as $t \rightarrow \infty$. For pure imaginary times, say $z = i\theta$ with $\theta \in \mathbb{R}$, we get the *Euler’s formula*

$$e^{i\theta} = \left(1 - \frac{\theta^2}{2!} + \frac{\theta^4}{4!} - \dots\right) + i\left(\theta - \frac{\theta^3}{3!} + \frac{\theta^5}{5!} - \dots\right) = \cos(\theta) + i\sin(\theta)$$

(and of course you may take the last identity as the “definition” of the trigonometric functions!). So, $\theta \mapsto e^{i\theta}$ defines a periodic function with period 2π , sending the real line \mathbb{R} onto the unit circle $\mathbb{S} = \{z \in \mathbb{C} \text{ s.t. } |z| = 1\}$. There follows from the group property that

$$\exp(t + i\theta) = e^t (\cos(\theta) + i\sin(\theta)) .$$

Finally, the exponential \exp is a periodic entire function with period $i2\pi$ which only omits the value 0, a holomorphic bijection of the cylinder $\mathbb{C}/i2\pi\mathbb{Z}$ onto $\mathbb{C} \setminus \{0\}$.

Interest rates and the exponential. Let x be the annual interest paid for a deposit (so that an interest of 0.2% mean $x = 0.02$). If the interest is paid once each year, an initial deposit of a euros increases to

$$a + xa = a \cdot (1 + x)$$

after one year. If, however, the interest is “computed” every six months, the same initial deposit produces

$$a + \frac{x}{2}a + \left(a + \frac{x}{2}a\right) \frac{x}{2} = a \cdot \left(1 + \frac{x}{2}\right)^2$$

after one year. By induction, we see that if the interest is computed every $12/n$ months, after one year we get a final capital of

$$a \cdot \left(1 + \frac{x}{n}\right)^n$$

The limit of the gain factor as $n \rightarrow \infty$,

$$E(x) = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$$

is another definition of the exponential function. If the argument lives in the Riemann sphere, you may think that $\exp(z) = (1 - z/\infty)^\infty$ has a zero of order ∞ at the point $p = \infty \in \mathbb{C}$.

Population dynamics. The exponential models the dynamics of a population in a unlimited environment. The *Malthusian/exponential model*¹⁷ is

$$\dot{N} = \lambda N$$

where $N(t)$ is the population at time t , and $\lambda > 0$ is some growth constant (the difference $\alpha - \beta$ between the natality rate and the mortality rate). The solution is $N(t) = N(0)e^{\lambda t}$. If we retire specimen at fixed rate $\alpha > 0$

$$\dot{N} = \lambda N - \alpha$$

we have a non-trivial stationary solution $\bar{N} = \alpha/\lambda$, and the difference $x(t) = N(t) - \bar{N}$ is still exponential.

This behaviour has to be compared with the *super-exponential model*

$$\dot{N} = \lambda N^2.$$

which undergoes a catastrophe (infinite population) in finite time! Indeed, the solution with $N(0) = N_0 > 0$ is $N(t) = N_0/(1 - \lambda t/N_0)$.

A more realistic model of population dynamics in a finite environment is the *logistic equation*¹⁸

$$\dot{N} = \lambda N(1 - N/M)$$

¹⁷T.R. Malthus, *An Essay on the Principle of Population*, London, 1798.

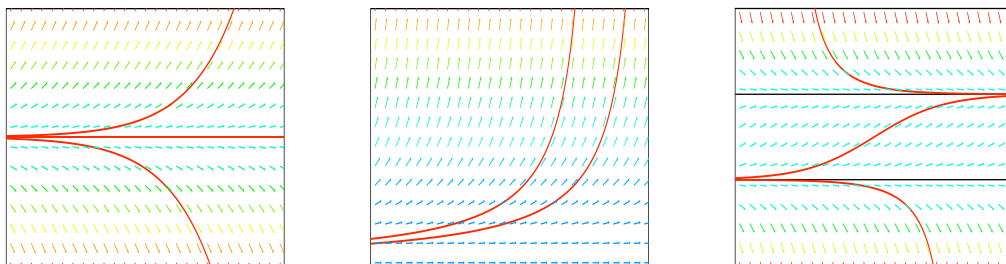
¹⁸P.F. Verhulst, Notice sur la loi que la population poursuit dans son accroissement, *Correspondance mathématique et physique* **10** (1838), 113-121.

where $\lambda > 0$ and the constant $M > 0$ is a maximal population. Observe that $\dot{N} \simeq \lambda N$ if $N \ll M$, and that $\dot{N} \rightarrow 0$ when $N \rightarrow M$. The relative population $x(t) = N(t)/M$ satisfies the “adimensional” logistic equation

$$\dot{x} = \lambda x(1 - x).$$

Here we see two equilibria: the trivial equilibrium $x(t) = 0$ and the maximum allowed population $x(t) = 1$. The generic solution with initial condition $0 < x(0) < 1$ is

$$x(t) = \frac{1}{1 + \left(\frac{1}{x_0} - 1\right) e^{-\lambda t}},$$



Exponential growth, super-exponential growth and logistic model.

2.4 Existence and uniqueness theorems

Solutions of a differential equation. Here we consider a generic first order ODE of the form

$$\dot{x} = v(x, t)$$

where the velocity field v is a (continuous) function defined in some extended phase space $X \times \mathbb{R}$. The phase space X may be some interval of the real line, an open subset of some Euclidean \mathbb{R}^n , or a differentiable manifold.

The problem we address is the existence and uniqueness of solutions of the initial value (or Cauchy) problem. A *local solution* passing through the point $(x_0, t_0) \in X \times \mathbb{R}$ is a solution $t \mapsto \varphi(t)$, defined in some neighbourhood I of t_0 , such that $\varphi(t_0) = x_0$. Eventually, we’ll be interested also in the possibility of extending such local solutions to larger intervals of times.

The basic existence theorem is ¹⁹

Theorem 2.2 (Peano). *Let $v(x, t)$ be a continuous velocity field in some domain A of the extended phase space \mathbb{R}^2 . Then for any point $(x_0, t_0) \in A$ passes at least one integral curve of the differential equation $\dot{x} = v(x, t)$.*

Proof. (Idea) Natural guesses for the solutions are Euler lines starting through (x_0, t_0) . If we restrict to a sufficiently small neighbourhood of (x_0, t_0) , we can assume that the velocity field is bounded, say $|v(x, t)| \leq K$, and that all such Euler lines lies in the “papillon” made of two triangles touching at (x_0, t_0) with slopes $\pm K$. Construct a family of Euler lines, graphs of $\varphi_n(t)$, such that the maximal step ε_n of the n -th line goes to 0 as $n \rightarrow \infty$. One easily sees that the family (φ_n) is bounded and equicontinuous. By the Ascoli-Arzelà theorem it admits a (uniformly) convergent subsequence. Finally, we claim that the sublimit $\varphi_{n_i} \rightarrow \varphi$ solves the differential equation. \square

¹⁹G. Peano, Sull’integrabilità delle equazioni differenziali del primo ordine, *Atti Accad. Sci. Torino* **21** (1886), 677-685. G. Peano, Demonstration de l’intégrabilité des équations différentielles ordinaires, *Mathematische Annalen* **37** (1890) 182-228.

Both existence and uniqueness may fail. The Hamilton-Jacobi equation

$$(\dot{x})^2 - xt + 1 = 0$$

cannot have solutions satisfying the initial condition $x(0) = 0$, for otherwise we would have a negative “kinetic energy” $(\dot{x})^2 = -1$ at that point!

Some regularity of the functions involved in a differential equation is also needed to ensure the uniqueness of solutions. For example, both curves $t \mapsto 0$ and $t \mapsto t^3$ solve the equation

$$\dot{x} = 3x^{2/3}$$

with initial condition $x(0) = 0$. The problem here is that the velocity field $v(t, x) = 3x^{2/3}$, although continuous, is not differentiable and not even Lipschitz at the origin.

Uniqueness of solutions. A velocity field $v(t, x)$, defined in a domain $I \times D$ of the extended phase space $\mathbb{R} \times \mathbb{R}^n$, is *locally Lipschitz* w.r.t. to the variable x if for any $(t_0, x_0) \in I \times D$ there is a neighbourhood $J \times U \ni (t_0, x_0)$ and a constant $L \geq 0$ such that

$$\|v(t, x) - v(t, y)\| \leq L \cdot \|x - y\| \quad \forall (t, x), (t, y) \in J \times U$$

If $v(t, x)$ has continuous derivative w.r.t. x , i.e. if the Jacobian

$$D_x v(t, x) = \left(\frac{\partial v_i}{\partial x_j}(t, x) \right)$$

exists and is continuous, then $v(t, x)$ is locally Lipschitz in any compact convex domain $I \times K \subset \mathbb{R} \times \mathbb{R}^n$. The basic uniqueness theorem is the following classical result by Lindelöf²⁰ and Picard.

Theorem 2.3 (Picard-Lindelöf). *Let $v(t, x)$ be a continuous velocity field defined in some domain D of the extended phase space $\mathbb{R} \times X$. If v is locally Lipschitz (for example continuously differentiable) w.r.t. the second variable x , then there exist one and only one local solution of $\dot{x} = v(t, x)$ passing through any point $(t_0, x_0) \in D$.*

Geometrically, the uniqueness theorem says that through any point (t_0, x_0) of the domain D there pass one and only one solution. Hence solutions, considered as curves in the extended phase space, cannot intersect each other.

In a domain where Picard’s theorem applies, if two local solutions agree in a common interval of times then they are indeed restrictions of a unique solution defined in the union of the respective domains. There follows that solutions are always extendible to a maximum domain. Such solutions are called *maximal solutions*.

Strategy of the proof of the Picard’s theorem. The first observation is that a function $\varphi(t)$ is a solution of the Cauchy problem for $\dot{x} = v(t, x)$ with initial condition $\varphi(t_0) = x_0$ iff

$$\varphi(t) = x_0 + \int_{t_0}^t v(s, \varphi(s)) ds$$

Now, we notice that the above identity is equivalent to the statement that φ is a fixed point of the so called *Picard’s map* $\phi \mapsto \mathcal{P}\phi$, sending a function $t \mapsto \phi(t)$ into the function

$$(\mathcal{P}\phi)(t) = x_0 + \int_{t_0}^t v(s, \phi(s)) ds$$

At this point, one must choose cleverly the domain of the Picard’s map, which is the space of functions where we think a solution should be. It will be a certain space \mathcal{C} of continuous functions, defined in an appropriate neighbourhood I of t_0 , equipped with a norm that makes it a complete metric space (hence a Banach space). The Lipschitz condition, together with continuity, satisfied by the velocity field will imply that if the interval I is sufficiently small then the Picard’s map $\mathcal{P} : \mathcal{C} \rightarrow \mathcal{C}$ is a contraction. The contraction principle (theorem 6.4) finally guarantees the existence and uniqueness of the fixed point of \mathcal{P} in \mathcal{C} .

²⁰M.E. Lindelöf, Sur l’application de la méthode des approximations successives aux équations différentielles ordinaires du premier ordre, *Comptes rendus hebdomadaires des séances de l’Académie des sciences* **114** (1894), 454-457. Digitized version online via <http://gallica.bnf.fr/ark:/12148/bpt6k3074>

Picard's iterations. The contraction principle actually says that the fixed point, i.e. the solution we are looking for, is the limit of any sequence $\phi, \mathcal{P}\phi, \dots, \mathcal{P}^n\phi, \dots$ of iterates of the Picard map starting with any initial guess $\phi \in \mathcal{C}$. In other words, the existence part of the theorem is “constructive”, it gives us a procedure to find out the solution, or at least a sequence of functions which approximate the solution.

Picard's iterations for simple ODEs. Consider the simple ODE $\dot{x} = v(t)$ with initial condition $x(t_0) = x_0$. Picard's recipe, starting from the initial guess $\phi(t) = x_0$ gives, already at the first step,

$$(\mathcal{P}\phi)(t) = x_0 + \int_{t_0}^t v(s)ds$$

which is the solution we know.

Picard's iterations for the exponential. Suppose you want to solve $\dot{x} = x$ with initial condition $x(0) = 1$. You start with the guess $\phi(t) = 1$, and then compute

$$(\mathcal{P}\phi)(t) = 1 + t \quad (\mathcal{P}^2\phi)(t) = 1 + t + \frac{1}{2}t^2 \quad \dots \quad (\mathcal{P}^n\phi)(t) = 1 + t + \frac{1}{2}t^2 + \dots + \frac{1}{n!}t^n$$

Hence the sequence converges (uniformly on bounded intervals) to the Taylor series of the exponential function

$$(\mathcal{P}^n\phi)(t) \rightarrow 1 + t + \frac{1}{2}t^2 + \dots + \frac{1}{n!}t^n + \dots = e^t,$$

which is the solution we already knew.

Details of the proof of the Picard's theorem. Choose a sufficiently small rectangular neighbourhood

$$I \times B = [t_0 - \varepsilon, t_0 + \varepsilon] \times \overline{B}_\delta(x_0)$$

around (t_0, x_0) , where $B = \overline{B}_\delta(x_0)$ denotes the closed ball with center x_0 and radius δ in X . There follows from continuity of v that there exists K such that

$$|v(t, x)| \leq K$$

for any $(t, x) \in I \times B$. There follows from the local Lipschitz condition for v that there exists M such that

$$|v(t, x) - v(t, y)| \leq M|x - y|$$

for any $t \in I$ and any $x, y \in B$. Now restrict, if needed, the (radius of the) interval I in such a way to get both the inequalities $K\varepsilon \leq \delta$ and $M\varepsilon < 1$. Let \mathcal{C} be the space of continuous functions $t \mapsto \phi(t)$ sending I into B . Equipped with the sup norm

$$\|\phi - \varphi\| = \sup_{t \in I} |\phi(t) - \varphi(t)|$$

this is a complete space. One verifies that the Picard's map sends \mathcal{C} into \mathcal{C} , since

$$|(\mathcal{P}\phi)(t) - x_0| \leq \int_{t_0}^t |v(s, \phi(s))| ds \leq K\varepsilon \leq \delta.$$

Finally, given two functions $\phi, \varphi \in \mathcal{C}$, one sees that

$$|(\mathcal{P}\phi)(t) - (\mathcal{P}\varphi)(t)| \leq \int_{t_0}^t |v(s, \phi(s)) - v(s, \varphi(s))| ds \leq M\varepsilon \sup_{t \in I} |\phi(t) - \varphi(t)|$$

hence $\|\mathcal{P}\phi - \mathcal{P}\varphi\| < M\varepsilon \|\phi - \varphi\|$. Since $M\varepsilon < 1$, this proves that the Picard's map is a contraction and the fixed point theorem allows to conclude. \square

We may not be able to solve them! Last but not least, we must keep in mind that we are not able to solve all equations. Actually, although we may prove the existence and the uniqueness for large classes of equations, we are simply not able to explicitly integrate the really interesting differential equations...

Ultimately we must recur to numerical methods to find approximate solutions and to qualitative analysis

Dependence on initial data and parameters Consider a family of ODEs

$$\dot{x} = v(t, x, \lambda)$$

where λ is a real parameter. We want to understand how solutions depend on the parameter λ . A basic instrument is the²¹

Theorem 2.4 (Grönwall's lemma). *Let $\phi(t)$ and $\psi(t)$ be two non-negative real valued functions defined in interval $[a, b]$ such that*

$$\phi(t) \leq K + \int_a^t \psi(s)\phi(s)ds$$

for any $a \leq t \leq b$ and some constant $K \geq 0$. Then

$$\phi(t) \leq K e^{\int_a^t \psi(s)ds}.$$

Proof. First, assume $K > 0$. Define

$$\Phi(t) = K + \int_a^t \psi(s)\phi(s)ds$$

and observe that $\Phi(a) = K > 0$, hence $\Phi(t) > 0$ for all $a \leq t \leq b$. The logarithmic derivative is

$$\frac{d}{dt} \log \Phi(t) = \frac{\psi(t)\phi(t)}{\Phi(t)} \leq \psi(t)$$

where we used the hypothesis $\phi(t) \leq \Phi(t)$. Integrating the inequality we get, for $a \leq t \leq b$,

$$\log \Phi(t) \leq \log \Phi(a) + \int_a^t \psi(s)ds.$$

Exponentiation gives the result, since

$$\phi(t) \leq \Phi(t) \leq K \cdot e^{\int_a^t \psi(s)ds}$$

The case $K = 0$ follows taking the limit of the above inequalities for a sequence of $K_n > 0$ decreasing to zero. \square

Continuous dependence on initial conditions. If $x(t)$ and $y(t)$ are two solutions of the same differential equation

$$\dot{x} = v(t, x)$$

then

$$x(t) - y(t) = x(0) - y(0) + \int_{t_0}^t (v(s, x(s)) - v(s, y(s))) ds$$

²¹T.H. Gronwall, Note on the derivative with respect to a parameter of the solutions of a system of differential equations, *Ann. of Math* **20** (1919), 292-296.

If $L(s)$ denotes the Lipschitz constant of $v(s, \cdot)$, we get

$$\|x(t) - y(t)\| \leq \|x(0) - y(0)\| + \int_{t_0}^t L(s) \|x(s) - y(s)\| ds$$

The Gronwall's lemma 2.4 gives the estimate

$$\|x(t) - y(t)\| \leq e^{\int_{t_0}^t L(s) ds} \|x(0) - y(0)\|$$

Observe that the above control also gives an alternative proof of uniqueness of solutions given a Lipschitz condition on the vector field.

Theorem 2.5 (smooth dependence on parameters). *Let $v(t, x, \lambda)$ be a family of vector fields defined on some domain of the extended phase space $D \subset \mathbb{R} \times X$ depending on a parameter $\lambda \in \Lambda \subset \mathbb{R}$. If v is of class C^k with $k \geq 1$, then in some neighbourhood of any $(t_0, x_0, \lambda_0) \in D \times \Lambda$ the local solutions of*

$$\dot{x} = v(t, x, \lambda)$$

with initial condition $x(t_0) = x_0$ are differentiable (indeed C^k) functions of (t, x, λ) .

A proof may be found in [BN05].

Warning. Continuous dependence does not exclude sensitive dependence on both initial conditions and parameters, even in the linear case! For example, the distance between solutions of $\dot{x} = \mu x$ with different $x(0)$ and/or μ may diverge for large time ...

2.5 Oscillations and cycles

The first remarkable natural phenomena are, of course, periodic motions.

Harmonic oscillator. The *harmonic oscillator* is the (phenomenon modeled by the) Newton equation

$$\ddot{q} = -\omega^2 q.$$

This is a quite universal equation, since it describes small oscillations around a “generic” stable equilibrium of any one-dimensional Newtonian system²² (indeed, take a Newton equation $m\ddot{x} = -dU'(x)$ of a particle in a potential field U . An equilibrium position is a zero of the force, i.e. a point x_0 where $U'(x_0) = 0$. It is “stable” if x_0 is a local minimum of the potential, so that the Taylor expansion of the potential around x_0 in powers of $q = x - x_0$ starts with $U(x) = \alpha + \frac{1}{2}\beta q^2 + \dots$, for some positive second derivative $U''(x_0) = \beta$. If we are only interested in small displacements of x around x_0 , we can safely disregard high order terms and approximate the Newton equation as $m\ddot{q} \simeq -\beta q$, which is an harmonic oscillator with resonant frequency $\omega = \sqrt{\beta/m}$).

Call $p = \dot{q}$ the momentum. The Newton equation $\ddot{q} = -\omega^2 q$ is equivalent to Hamilton's first order equations

$$\begin{aligned} \dot{q} &= p \\ \dot{p} &= -\omega^2 q. \end{aligned}$$

If we define the complex variable $z = \omega q + i\dot{q}$, Newton equation then takes the form of a first order linear equation in the complex line, namely $\dot{z} = -i\omega z$, whose solution is $z(t) = e^{-i\omega t} z(0)$.

²² “The harmonic oscillator, which we are about to study, has close analogs in many other fields; although we start with a mechanical example of a weight on a spring, or a pendulum with a small swing, or certain other mechanical devices, we are really studying a certain *differential equation*. This equation appears again and again in physics and other sciences, and in fact is a part of so many phenomena that its close study is well worth our while. Some of the phenomena involving this equation are the oscillations of a mass on a spring; the oscillations of charge flowing back and forth in an electrical circuit; the vibrations of a tuning fork which is generating sound waves; the analogous vibrations of the electrons in an atom, which generate light waves; the equations for the operation of a servosystem, such as a thermostat trying to adjust a temperature; complicated interactions in chemical reactions; the growth of a colony of bacteria in interaction with the food supply and the poison the bacteria produce; foxes eating rabbits eating grass, and so on; ...”

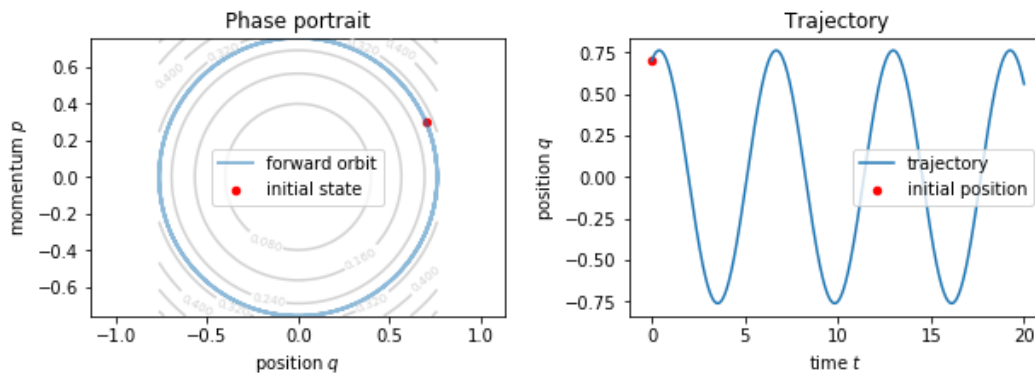
R.P. Feynman [Fe63]

In terms of the original (physical) variables, the solutions read

$$q(t) = q_0 \cos(\omega t) + \frac{v_0}{\omega} \sin(\omega t) = A \sin(\omega t + \phi)$$

where the amplitude A and the initial phase ϕ depend on the initial conditions $q(0) = q_0$ and $\dot{q}(0) = v_0$. So, all trajectories are periodic with common period $2\pi/\omega$, and orbits are ellipses in the q - \dot{q} plane, determined by the conserved energy

$$E = \frac{1}{2} (\dot{q}^2 + \omega^2 q^2) = \omega^2 A^2.$$



Harmonic oscillator, orbit and time series.

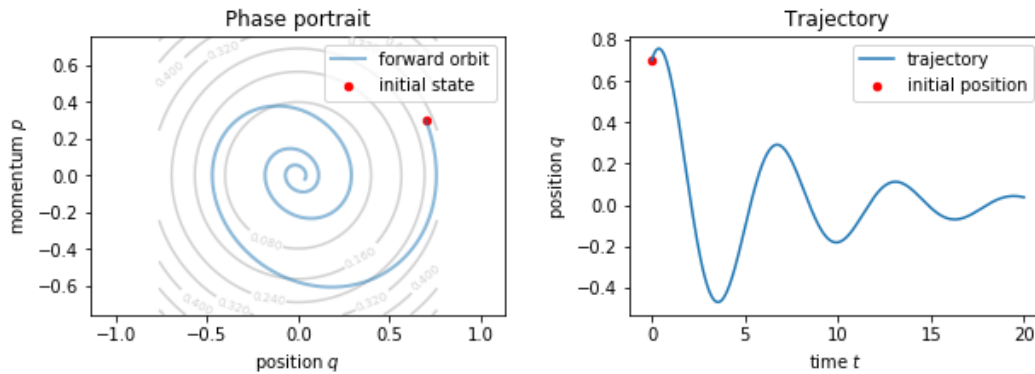
Damped oscillations. Adding friction to an harmonic oscillator we get

$$\ddot{q} = -2\alpha\dot{q} - \omega^2 q,$$

where $\alpha > 0$ is some friction coefficient. The guess $q(t) = e^{-\alpha t} y(t)$ gives $\ddot{y} = \delta y$ where the “discriminant” is $\delta = \omega^2 - \alpha$. Find the general solution, draw pictures and discuss the cases $\alpha^2 < \omega^2$ (under-critical damping), $\alpha^2 = \omega^2$ (critical damping), and $\alpha^2 > \omega^2$ (overcritical damping). Show that the energy

$$E(q, \dot{q}) = \frac{1}{2} \dot{q}^2 + \frac{1}{2} \omega^2 q^2$$

decreases with time outside equilibrium points.



Underdamped harmonic oscillator, orbit and time series.

Mathematical pendulum. The Newton equation

$$I\ddot{\theta} = -mg\ell \sin \theta$$

models the motion of an idealized pendulum (meaning a point mass attached to a wire of negligible weight, under a constant gravitational force) with mass m and length ℓ , where $I = m\ell^2$ is the moment of inertia, g is the gravitational acceleration (near the Earth's surface), and θ is the angle of the wire with the origin $\theta = 0$ located at the stable equilibrium point. The energy

$$E = \frac{1}{2}\dot{\theta}^2 - mg\ell \cos \theta$$

is a constant of the motion. We can define the resonant frequency $\omega = \sqrt{mg\ell/I} = \sqrt{g/\ell}$ and write the equation as

$$\ddot{\theta} = -\omega^2 \sin \theta$$

Observe that in the limit of small oscillations we could replace $\sin \theta \simeq \theta$ and we are back to the harmonic oscillator $\ddot{\theta} = -\omega^2 \theta$. To simplify things, let's take $\omega = 1$. Solving the energy for $\dot{\theta}^2$ we see that the motion with energy E is given implicitly by the "elliptic integral"

$$t = \int \frac{d\theta}{\sqrt{2(E - \cos(\theta))}}$$

What does a mathematician/physicist do when he/she face an integral and doesn't see how to solve in terms of known functions? He/she gives a name to it.

Define $k = \sqrt{\frac{E+1}{2}}$ and then $x = \frac{1}{k} \sin(\theta/2)$. The conservation of energy reads

$$\dot{x} = \sqrt{(1-x^2)(1-k^2x^2)}$$

There follows that time is given by the so called *Jacobi's elliptic integral of the first kind*

$$t = \int \frac{dx}{\sqrt{(1-x^2)(1-k^2x^2)}}$$

The solution, actually the inverse function $x = \text{sn}(t, k)$ as a function of t and the parameter k , is "named" *Jacobi elliptic function*.

This is the beginning of a long and interesting story. You may want to know that sn , as well its relatives, is a quotient of products of *Jacobi's theta functions*, hence, we are at the intersection between complex analysis, algebraic geometry, number theory, ...

Kepler problem. Kepler problem deals with the motion of two point-like bodies (planets and/or stars) under mutual gravitational interaction. Let $m_1, m_2 > 0$ be their masses, and $q_1, q_2 \in \mathbb{R}^3$ their positions, respectively. Gravitational interaction is described by the conservative force $-\nabla V$ with potential energy

$$V(q_1, q_2) = G \frac{m_1 m_2}{|q_1 - q_2|}$$

where G is the gravitational constant. This force verifies the "third law of dynamics", hence the total linear and angular momentum

$$P = m_1 \dot{q}_1 + m_2 \dot{q}_2 \quad \text{and} \quad M = m_1 q_1 \wedge \dot{q}_1 + m_2 q_2 \wedge \dot{q}_2$$

are conserved. This implies that the center of mass moves at uniform rectilinear speed and that the motion of the two bodies takes place in a plane orthogonal to the angular momentum M . If we choose a Galileian reference system where $P = 0$ and M is parallel to the z -axis (in particular M is supposed different from the zero vector, a case which leads to a collision ...), the full system is described by the single vector $q_2 - q_1$ in the x - y plane, which we write in polar coordinates as $\rho e^{i2\pi\theta}$. It turns out that the two-body problem is equivalent to the motion of a single point mass $m = \frac{m_1 m_2}{m_1 + m_2}$ moving on a plane under the influence of a potential energy $V(\rho) = -G \frac{m}{\rho}$, the (conserved) energy being

$$E = \frac{1}{2}m(\dot{\rho}^2 + \rho^2 \dot{\theta}^2) + V(\rho)$$

Observe that if one of the bodies is much bigger than the other (like the Sun and the Earth), say $m_1 \gg m_2$, then the center of mass nearly coincides with the position q_1 of the bigger body, while the reduced mass m is essentially the mass m_2 of the smaller one (hence it looks like the Earth moving around the Sun, as Galileo had suggested).

Central forces. Consider the Newton equation

$$m\ddot{\mathbf{r}} = F(|\mathbf{r}|) \hat{\mathbf{r}}$$

describing the motion of a particle (planet) of mass m in a central force field F . Conservation of angular momentum implies that the motion is planar, hence we may take $\mathbf{r} \in \mathbb{R}^2$. In polar coordinates $\mathbf{r} = \rho e^{i\theta}$, the equations read

$$\begin{aligned} \ddot{\rho} - \rho\dot{\theta}^2 &= F(\rho)/m \\ \rho\ddot{\theta} + 2\dot{\rho}\dot{\theta} &= 0. \end{aligned}$$

The second equation says that the “areal velocity” $\ell = \rho^2\dot{\theta}$ is a constant of the motion (Kepler’s second law).

Taking Newton’s gravitational force $F(\rho) = -\frac{GmM}{\rho^2}$ (where M is the mass of the Sun and G is the gravitational constant), the first equation may be written as

$$m\ddot{\rho} = -\frac{\partial}{\partial \rho} V_\ell(\rho),$$

where we defined the “effective potential energy” as $V_\ell(\rho) = \frac{1}{2}m\frac{\ell^2}{\rho^2} - G\frac{mM}{\rho}$. The conserved energy is

$$E = \frac{1}{2}m\dot{\rho}^2 + \frac{1}{2}m\frac{\ell^2}{\rho^2} - G\frac{mM}{\rho}.$$

Now we set $\rho = 1/x$ and look for a differential equation for x as a function of θ . Computation shows that $dx/d\theta = -\dot{\rho}/\ell$, and, using conservation of ℓ , that $d^2x/d\theta^2 = -\rho^2\ddot{\rho}/\ell^2$. There follows that the first Newton equation reads

$$\frac{d^2x}{d\theta^2} + x = -\frac{1}{\ell^2 x^2 m} F(1/x).$$

we get

$$\frac{d^2x}{d\theta^2} + x = -\frac{GM}{\ell^2}.$$

The general solution of this second order linear differential equation is

$$x(\theta) = \frac{GM}{\ell^2} (1 + e \cos(\theta - \theta_0)),$$

for some constants e and θ_0 . Back to the original radial variable we get the solution

$$\rho(\theta) = \frac{\ell^2/GM}{1 + e \cos(\theta - \theta_0)},$$

Hence, orbits are conic sections with eccentricity e and focus at the origin: an ellipse for $0 \leq e < 1$ (corresponding to negative energy, hence to planets, and this is Kepler’s first law), a parabola for $e = 1$ (corresponding to zero energy), an hyperbola for $e > 1$ (corresponding to positive energy).

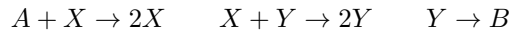
2.6 Phenomenological models

A number of phenomenological models (i.e. models which are not fundamental laws of nature), like the ones below, are also a source of interesting dynamical behaviour.

Lotka-Volterra predator-prey model. The *Lotka-Volterra system* is the first-order non-linear differential equation

$$\begin{aligned}\dot{x} &= ax - bxy \\ \dot{y} &= -cy + dxy\end{aligned}$$

It has been proposed by Vito Volterra²³ to model competition between x preys and y predators, and by Alfred J. Lotka²⁴ to model the cyclic behavior of certain chemical reactions, like the abstract scheme

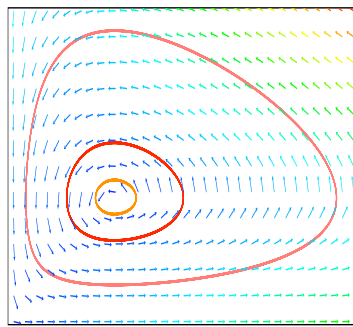


Preys increase exponentially at rate a and are killed at rate proportional to the probability of being captured by a predator, while predators decrease exponentially at rate c and increase at rate proportional to the probability of capturing preys.

The function

$$H(x, y) = dx + by - c \log x - a \log y$$

is a constant of the motion, i.e. $\frac{d}{dt}H(x(t), y(t)) = 0$. There follows that orbits are contained (and actually are) in the level curves of H .



Phase portrait of the Lotka-Volterra system.

ex: Discuss the possible dynamics depending on the values of the parameters.

Competing species. Competition between two species sharing the same environment could be modelled by a system of coupled logistic equations

$$\begin{aligned}\dot{x} &= \lambda x(1 - x) - \beta xy \\ \dot{y} &= \mu y(1 - y) - \gamma xy\end{aligned}$$

ex: Try to understand the phase portrait given some different values of the parameters.

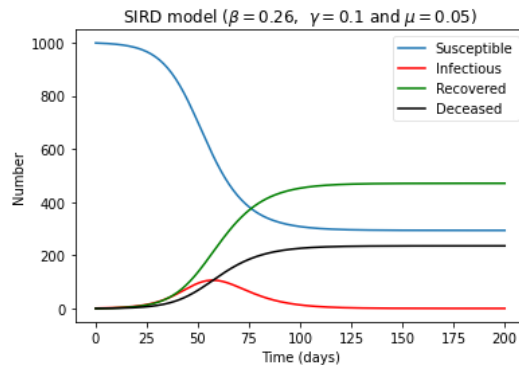
SIRD models. A simple model of an epidemic in a fixed population of N specimen is the so called *SIRD model*, the system

$$\begin{aligned}\dot{S} &= -\beta SI/N \\ \dot{I} &= \beta SI/N - \gamma I - \mu I \\ \dot{R} &= \gamma I \\ \dot{D} &= \mu I\end{aligned}$$

Here $S(t)$ is the susceptible population, $I(t)$ the infected population, $D(t)$ the deceased population, and β, γ and μ are convenient positive parameters.

²³V. Volterra, Variazioni e fluttuazioni del numero d'individui in specie di animali conviventi, *Mem. Acad. Lincei* **2** (1926), 31-113. V. Volterra, *Leçons sur la Théorie Mathématique de la Lutte pour la Vie*, Paris 1931.

²⁴A.J. Lotka, *J. Amer. Chem. Soc.* **27** (1920), 1595. A.J. Lotka, *Elements of physical biology*, Williams & Wilkins Co. 1925.

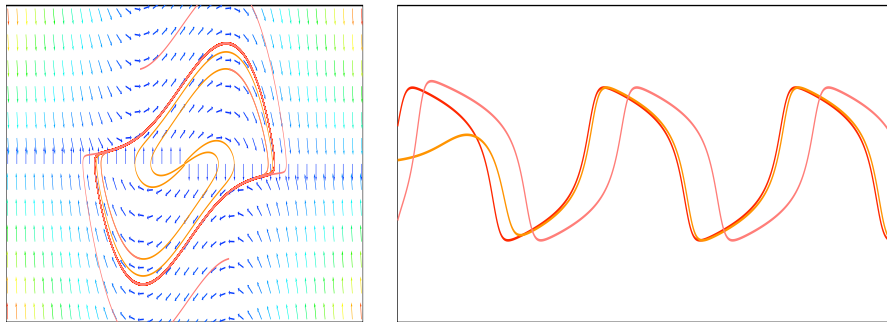


A trajectory of a SIRD model.

Van der Pol's oscillator. The *van der Pol oscillator*²⁵ is the second-order non-linear differential equation

$$\ddot{q} - \mu(1 - q^2)\dot{q} + q = 0$$

which models current in a circuit with a non-linear element.



Phase portrait and time series of the Van der Pol's oscillator.

Brusselator. The *Brusselator* is an auto-catalytic model proposed by Ilya Prigogine and collaborator²⁶ which models the abstract reaction



and reads

$$\begin{aligned} \dot{x} &= \alpha - (\beta + 1)x + x^2y \\ \dot{y} &= \beta x - x^2y \end{aligned}$$

ex: Observe what happens to the concentrations X e Y , namely x and y , when the concentrations $[A] \sim \alpha$ and $[B] \sim \beta$ are kept constant.

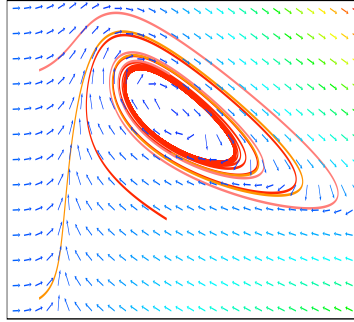
ex: Simulate the system

$$\begin{aligned} \dot{x} &= \alpha - (b + 1)x + x^2y \\ \dot{y} &= bx - x^2y \\ \dot{b} &= -bx + \delta \end{aligned}$$

for the concentrations of X , Y and B , obtained when the concentration $[A] \sim \alpha$ is maintained constant and B is injected with constant velocity $v \sim \delta$.

²⁵B. van der Pol, A theory of the amplitude of free and forced triode vibrations, *Radio Review* **1** (1920), 701-710 and 754-762. B. van der Pol and J. van der Mark, Frequency demultiplication, *Nature* **120** (1927), 363-364.

²⁶I. Prigogine and R. Lefever, Symmetry breaking instabilities in dissipative systems, *J. Chem. Phys.* **48** (1968), 1655-1700. P. Glansdorff and I. Prigogine, *Thermodynamic theory of structure, stability and fluctuations*, Wiley, New York 1971. G. Nicolis and I. Prigogine, *Self-organization in non-equilibrium chemical systems*, Wiley, New York 1977.

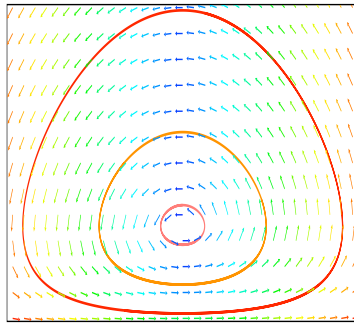


Phase portrait of the Brussellator.

Goodwin oscillator. A system modeling the interaction protein-mRNA was proposed by Goodwin²⁷

$$\begin{aligned}\dot{M} &= \frac{1}{1+P} - \alpha \\ \dot{P} &= M - \beta\end{aligned}$$

where M and P denote the relative concentrations of mRNA and protein, respectively.

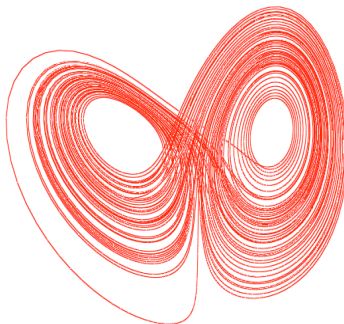


Phase portrait of the Goodwin oscillator.

Lorenz attractor. Finally, we mention the *Lorenz system*²⁸

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= x(\rho - z) - y \\ \dot{z} &= xy - \beta z\end{aligned}$$

For values of the parameters like $\sigma \simeq 10$, $\rho \simeq 28$ and $\beta \simeq 8/3$, one observe trajectories which diverge from one another, and yet oscillate all along the figure-eight above.



²⁷B.C. Goodwin, *Temporal organization in cells*, Academic Press, London/New York 1963. B.C. Goodwin, Oscillatory behaviour in enzymatic control processes, *Adv. Enzyme Regul.* **3** (1965), 425-438.

²⁸E.N. Lorenz, Deterministic nonperiodic flow, *J. Atmospheric Science* **20** (1963), 130-141.

Some orbits of the Lorenz attractor.

This strange phenomenon motivated an important part of the modern theory of dynamical systems.

3 Topological dynamical systems, basic definitions

3.1 Transformations

Transformations. In these notes, we'll be mainly interested in discrete time dynamical systems, i.e. actions of \mathbb{N}_0 or \mathbb{Z} on some space X , generated by a transformation/map

$$f : X \rightarrow X.$$

Apart from some special cases, X will be a topological space (or even a metric space). The transformation will be continuous, or at least piecewise continuous. In such cases we speak of *topological dynamical system*.

The “(forward) iterates” of a transformation f are the transformations $f^n : X \rightarrow X$, with $n \in \mathbb{N}_0$, defined inductively according to

$$f^0 = \text{id} \quad \text{and} \quad f^{n+1} = f \circ f^n \quad \text{if } n \geq 0$$

(warning! with this notation $f^2(x)$ is not the square of $f(x)$, but $f(f(x)) \dots$).

In general, if $n \in \mathbb{N}$ and $A \subset X$, then $f^{-n}(A)$ denotes the set

$$f^{-n}(A) = \{x \in X \text{ s.t. } f^n(x) \in A\}.$$

If f is invertible (e.g. is an homeomorphism), we can also define the backward iterates, and therefore the transformations $f^n : X \rightarrow X$ for all $n \in \mathbb{Z}$.

We have therefore an action $\Phi : \mathbb{N}_0 \times X \rightarrow X$, or $\Phi : \mathbb{Z} \times X \rightarrow X$ if f is invertible, defined by $\Phi_n(x) = f^n(x)$.

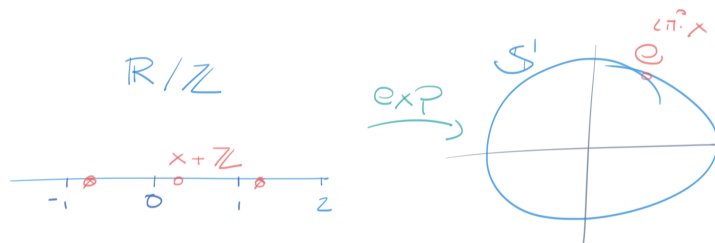
Phase/states space. In the following, (X, d) will be a metric space equipped with its natural topology τ , locally compact (any point admits a compact neighbourhood) and separable (admits a countable dense subset, and therefore, being a metric space, a countable basis for the topology). For example, regions of \mathbb{R}^N , intervals of the line, the circle \mathbb{R}/\mathbb{Z} , the torus $\mathbb{R}^N/\mathbb{Z}^n$, the complex plane \mathbb{C} , the Riemann sphere $\overline{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$, Cantor sets, and Cartesian products of finite spaces. Also, in order to avoid trivialities, we'll always assume tacitly that X is not a finite set.

Translations in homogeneous spaces. The simplest, tautological, way to build actions is algebraic. Let G be a topological group (a group equipped with a Hausdorff topology such that the group operations $(g, g') \mapsto gg'$ and $g \mapsto g^{-1}$ are continuous). Given a closed subgroup $\Gamma \subset G$, one can consider the homogeneous space $X = G/\Gamma = \{g\Gamma; g \in G\}$, equipped with the quotient topology (the finest topology in G/Γ such that the projection $\pi : G \rightarrow G/\Gamma$ is continuous). If Γ is not too large or wild, for example if Γ is discrete, X is a sufficiently large and interesting space.

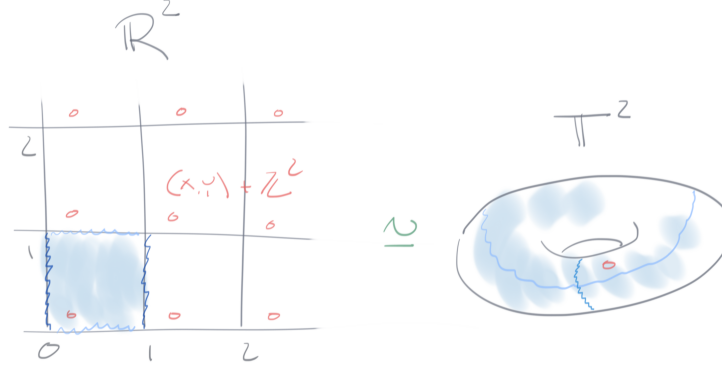
Every subgroup $S \subset G$ acts on the homogeneous space $X = G/\Gamma$, the action $S \times G/\Gamma \rightarrow G/\Gamma$ being $(s, g\Gamma) \mapsto sg\Gamma$. The space of orbits is the quotient $S \backslash G/\Gamma$.

In particular, a cyclic subgroup $S = \{s^n\}_{n \in \mathbb{Z}}$ generates an action $\Phi : \mathbb{Z} \times X \rightarrow X$ defined by $\Phi_n(g\Gamma) = s^n g\Gamma$, which consists in iterating the left translations $g\Gamma \mapsto sg\Gamma$ of a generator.

Translations of the torus. The N -dimensional *torus* is the quotient space $\mathbb{T}^N := \mathbb{R}^N/\mathbb{Z}^N$ of the additive Abelian group \mathbb{R}^N modulo the discrete subgroup \mathbb{Z}^N of integer vectors. Observe that the torus is itself an Abelian group. The one-dimensional case $\mathbb{T}^1 = \mathbb{R}/\mathbb{Z}$ is also called *circle*, because it is isomorphic to the unit circle of the complex plane, the multiplicative Abelian group $\mathbb{S}^1 = \{z \in \mathbb{C}, |z| = 1\}$, under the exponential map $x \mapsto e^{2\pi i x}$.



In dimension two, we may notice that any class $(x, y) + \mathbb{Z}^2$ has a unique representative in the unit square $[0, 1) \times [0, 1)$, and that opposite sides of the closed unit square must be identified according to $(0, y) \sim (1, y)$ and $(x, 0) \sim (x, 1)$. The resulting quotient space, the two-torus $\mathbb{T}^2 = \mathbb{R}^2 / \mathbb{Z}^2$, is therefore the surface of a donut.



Any $\alpha \in \mathbb{R}^N$ defines a translation $T_\alpha : \mathbb{R}^N \rightarrow \mathbb{R}^N$, according to $T_\alpha(x) = x + \alpha$. As explained above, the translation defines a *rotation* $R_\alpha : \mathbb{T}^N \rightarrow \mathbb{T}^N$, according to $R_\alpha(x + \mathbb{Z}^N) := x + \alpha + \mathbb{Z}^N$.

Generic properties. We will often want to talk about “most trajectories”, or “almost all trajectories”.

Being X a topological space, one could consider (probability or infinite) measures on the Borel σ -algebra of X . Given such a measure μ , one says that a properties is satisfied for μ -almost all points if the subset $N \subset X$ of those points which do not have the property has measure $\mu(N) = 0$.

The topological counterpart of the dichotomy “zero-one probability” is possible when X is a *Baire space*, i.e. a Hausdorff (any two distinct points have disjoint neighbourhoods) topological space where a countable intersection of dense open sets is dense. Baire theorem, that we state and prove for the reader’s convenience, says that examples of Baire spaces are complete metric spaces.

Theorem 3.1 (Baire). *Let X be a complete metric space. The intersection of a countable family of open and dense subsets is dense in X .*

Proof. Let A_n , with $n \in \mathbb{N}$, be open and dense subsets of X . Let U be any non-empty open subset of X . Since A_1 is dense and open, there exists a point $x_1 \in X$ and a positive radius $\varepsilon_1 < 1$ such that

$$\overline{B_{\varepsilon_1}(x_1)} \subset A_1 \cap U$$

Inductively, one shows that we may find a sequence of points x_n and positive radii $\varepsilon_n < 1/n$ such that

$$\overline{B_{\varepsilon_{n+1}}(x_{n+1})} \subset A_{n+1} \cap B_{\varepsilon_n}(x_n)$$

By the Cantor intersection theorem, there exists (at least) a point $x \in \bigcap_{n \geq 1} \overline{B_{\varepsilon_n}(x_n)}$, and by construction this point belongs to U as well as to the intersection of all the A_n ’s. \square

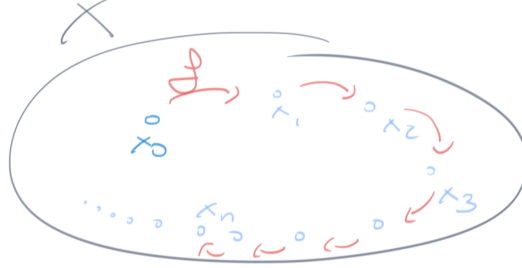
A subset $R \subset X$ is said *residual* if it contains a countable intersection of dense open sets. A subset $M \subset X$ is said *meager* if it is a countable union of “nowhere dense” subsets (subsets such that the closure has empty interior), i.e. if its complement $X \setminus M$ is residual. A property is said *generic* if the subset $P \subset X$ of those points with this property is residual.

3.2 Trajectories and orbits

Trajectories. Given a transformation $f : X \rightarrow X$, we are mainly interested in the asymptotic behavior of the “history” of a point $x \in X$, the sequence of points

$$x \mapsto f(x) \mapsto f^2(x) \mapsto f^3(x) \mapsto \dots$$

obtained recursively applying f to the point x . If X is the space state of a physical system, and if the system is (prepared) in the state x at time $t = 0$, then it will be in the state $f(x)$ at time 1, in the state $f^2(x) = f(f(x))$ at time 2, and so on.



The *trajectory* of $x \in X$ is the sequence $(x_n)_{n \in \mathbb{N}_0}$, the function that, given the “initial condition” $x_0 = x$, produces the states $x_n = f^n(x)$ of the system at each time $n \geq 0$. Thus, the trajectory of x is the solution of the recurrence

$$x_{n+1} = f(x_n)$$

with initial condition $x_0 = x$.

Orbits. The *forward/positive orbit* of $x \in X$ is the image of its trajectory, i.e. the set

$$O_f^+(x) := \{f^n(x)\}_{n \in \mathbb{N}_0}$$

(we put the superscript “+” to remind that we are only allowed to go forward in time, since in general f will not be invertible). It is the “future” of a point.

A point x may have more than one pre-image, and therefore its “past” is not unique. The *full orbit* of a point $x \in X$ is the set

$$O_f(x) := \{x' \in X : \exists n, m \geq 0 : f^n(x') = f^m(x)\}$$

i.e the set of points which have eventually the same future of x .

If f is invertible, the full orbit coincides with the *complete orbit* of a point x , defined as

$$O_f(x) := \{f^n(x)\}_{n \in \mathbb{Z}}$$

the past and future of a point.

Observe that “being in the same full-orbit” is an equivalence relation, and therefore X is a disjoint union of equivalence classes, i.e. orbits. It must be said that the quotient space, the space of orbits X/f , may be messy if trajectories are not regular (and this is when things get interesting!). For example, if there exists a dense orbit, then the quotient topology in X/f is the trivial topology. Thus, the space of orbits, as a topological space, does not contain much informations on the dynamics of the system.

3.3 Periodic orbits and basin of attraction

Fixed points. The simplest orbits are (composed of) *fixed points* of f , those states $p \in X$ such that

$$f(p) = p.$$

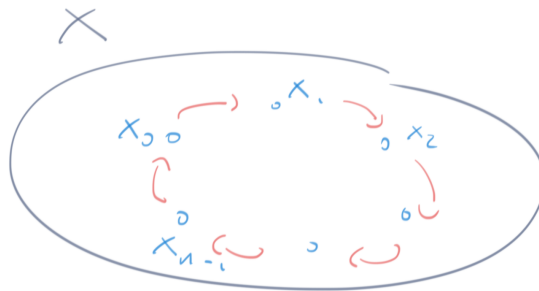
Geometrically, fixed points are the intersections of the graph of f with the “diagonal” $\Delta \subset X \times X$. If X is a linear space, fixed points are roots of the equation $f(x) - x = 0$.

The set of fixed points of f is denoted by $\text{Fix}(f) \subset X$. Since f is continuous, it is a closed subset of X .

Periodic orbits. A point $p \in X$ is said *periodic* if it is a fixed point of some iterate f^k , i.e. if it belongs to some $\text{Fix}(f^k)$. A periodic point p is periodic with *period* $n \geq 1$ if $f^n(p) = p$ and n is the smallest of those times $k \geq 1$ such that $f^k(p) = p$. Thus, the forward orbit of the periodic point p is a *cycle*, a finite set

$$\pi = O^+(p) = \{p, f(p), f^2(p), \dots, f^{n-1}(p)\}$$

of points which are permuted by the transformation f . The cardinality $|\pi| = n$ of the periodic orbit π is the common period of its points.



A point x may have a finite orbit without being periodic: this happens when there exists a time $k \geq 1$ such that $f^k(x)$ is a periodic point. Such points are called *pre-periodic*.

It is convenient to denote $\text{Per}_n(f) := \text{Fix}(f^n)$ the set of fixed points of the transformation f^n , called “ n -periodic points”, that is the set of those periodic points of f whose period divides n . Then

$$\text{Per}(f) = \bigcup_{n \geq 1} \text{Per}_n(f)$$

denotes the set of periodic points of the map f . Observe that any of the sets $\text{Per}_n(f)$ is closed, because f^n is continuous, but their union $\text{Per}(f)$ may not be closed.

It will be interesting, later, to compute or estimate the cardinalities $P_n(f) := \text{card}(\text{Per}_n(f))$, provided they are finite. Also interesting will be the cardinalities $\Pi_n(f)$ of periodic orbits π of length $|\pi| = n$. Clearly, $P_n(f) = \sum_{m|n} m \Pi_m(f)$.

Convergent trajectories. If a trajectory is convergent, then its limit is a fixed point of f . Indeed, if $f^n(x) \rightarrow p$, the continuity of f implies that

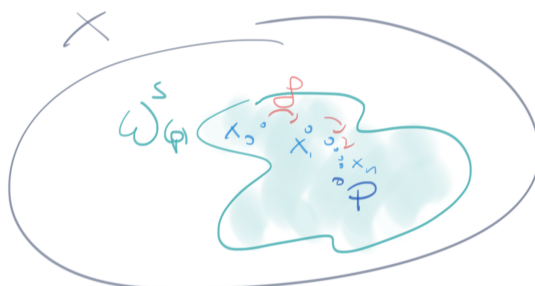
$$f(p) = f\left(\lim_{n \rightarrow \infty} f^n(x)\right) = \lim_{n \rightarrow \infty} f^{n+1}(x) = p.$$

Basin of attraction. Let p be a fixed point of $f : X \rightarrow X$. The *basin of attraction*, or *stable set*, of p is the set of those points $x \in X$ whose trajectories converge to p , i.e.

$$W^s(p) := \left\{ x \in X \text{ s.t. } \lim_{n \rightarrow \infty} f^n(x) = p \right\}$$

Uniqueness of limits of convergent sequences in a metric space implies that stable sets of different fixed points are disjoint.

In an obvious manner one defines the basin of attraction of a periodic orbit.



Endomorphisms of linear spaces. Between the simplest dynamical systems are endomorphisms of a linear space. For example, endomorphisms of \mathbb{R}^N are defined, in the canonical basis, by matrices $A \in \text{Mat}_{n \times n}(\mathbb{R})$, according to $f(x) = Ax$ (vectors are column vectors, and the product is the usual product between matrices). The origin is a fixed point, by linearity. Other fixed points are the eigenvectors with eigenvalue $\lambda = 1$, non-trivial solutions of the homogeneous equation $Ax = x$. Periodic points with period n are eigenvectors with eigenvalue λ such that $\lambda^n = 1$.

Collatz/Kakutani/Syracuse/Ulam problem. Consider the *Collatz map* $f : \mathbb{N} \rightarrow \mathbb{N}$, defined as

$$f(n) = \begin{cases} n/2 & \text{if } n \text{ is even} \\ 3n + 1 & \text{if } n \text{ is odd} \end{cases}$$

It is clear that $4 \mapsto 2 \mapsto 1$ is a cycle. *Collatz conjecture* (“... an extraordinarily difficult problem, completely out of reach of present day mathematics”, according to Lagarias²⁹) affirms that this is the only cycle and that any initial condition will eventually fall in this cycle.

Dynamics on a finite state spaces. Dynamics in a finite state space is almost trivial. Consider a transformation $f : X \rightarrow X$ of a finite set $X \approx \{1, 2, \dots, N\}$. It is clear than any trajectory is eventually periodic, so that the dynamics is completely described by a finite set of periodic orbits and their disjoint basins of attraction. The particular case of invertible transformations consists essentially in the study of the symmetric group S_N , permutations of X .

Nevertheless, if the map is chosen randomly, according to the natural uniform probability measure in the space X^X , we may compute mean values of the number of attractors, their size and sizes of their basins ... This is the *random map model*, a source of interesting mathematical and physical problems.

ex: (Affine maps) Draw some orbits of the transformations of the complex plane $f : \mathbb{C} \rightarrow \mathbb{C}$ defined by

$$f(z) = z + \alpha \quad \text{or} \quad f(z) = \lambda z$$

for different values of the parameters $\alpha, \lambda \in \mathbb{C}$. Explain for which values of those parameters there exists periodic orbits.

ex: Find, when possible, periodic orbits (of small period) of the transformations of the interval

$$f(x) = \pm x^3 \quad f(x) = x^{1/3} \quad f(x) = x^3 \pm x$$

$$f(x) = x^2 + 1/4 \quad f(x) = |1 - x| \quad f(x) = x^2 - 2 \quad f(x) = \sin x \quad f(x) = \cos x$$

$$f(x) = x(1 - x) \quad f(x) = 2x(1 - x) \quad f(x) = 3x(1 - x) \quad f(x) = 4x(1 - x)$$

ex: Find the basins of attraction of the fixed points of $f(x) = x^2$ and $f(x) = x^3$, considered as transformations of the real line \mathbb{R} .

ex: Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be the linear homogeneous transformation of the line defined by $f(x) = \lambda x$. Study the basin of attraction of $p = 0$ depending on the “multiplier” λ .

ex: Do the same for $f(z) = \lambda z$ defined in the complex line \mathbb{C} .

²⁹J.C. Lagarias (ed.), *The Ultimate Challenge: The $3x + 1$ Problem*, AMS, 2010.

ex: Find the basin of attraction of the origin for the linear maps $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined, in the canonical basis, by the following 2×2 real matrices:

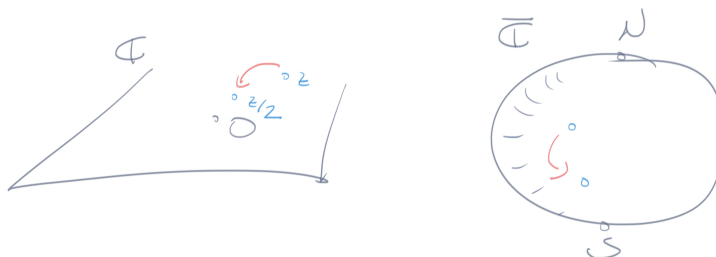
$$\begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix} \quad \begin{pmatrix} 2 & 0 \\ 0 & 1/3 \end{pmatrix} \quad \begin{pmatrix} 1/2 & 0 \\ 0 & 1/3 \end{pmatrix} \quad \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} \quad \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \quad \frac{1}{5} \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \quad 2 \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

North-South map. The map $z \mapsto z/2$ of the complex plane extends to an automorphism $f : \overline{\mathbb{C}} \rightarrow \overline{\mathbb{C}}$ of the Riemann sphere $\overline{\mathbb{C}} := \mathbb{C} \cup \{\infty\}$, declaring that $f(\infty) = \infty$. It is clear that the basin of attraction of 0 is all of $\mathbb{C} = \overline{\mathbb{C}} \setminus \{\infty\}$.

The stereographic projection $\pi : \mathbb{S}^2 \setminus \{N\} \rightarrow \mathbb{C}$ extends to a bijection between the two-sphere $\mathbb{S}^2 = \{x^2 + y^2 + z^2 = 1\} \subset \mathbb{R}^3$ and the Riemann sphere $\overline{\mathbb{C}}$, sending the North-pole $N = (0, 0, 1)$ to $\pi(N) = \infty$ and the South-pole $S = (0, 0, -1)$ to $\pi(S) = 0$. The composition $g := \pi^{-1} \circ f \circ \pi : \mathbb{S}^2 \rightarrow \mathbb{S}^2$ is called *North-South map*. It fixes the North-pole and the South-pole, and the orbit of any other point converges (along meridians) to the South-pole. Thus, the basin of attraction of the South-pole is $W^s(S) = \mathbb{S}^2 \setminus \{N\}$.

Sometimes, also the restriction of g to a meridian (for example, the meridian corresponding to the real line under stereographic projection), which is a self-map of the circle, is called North-South map.

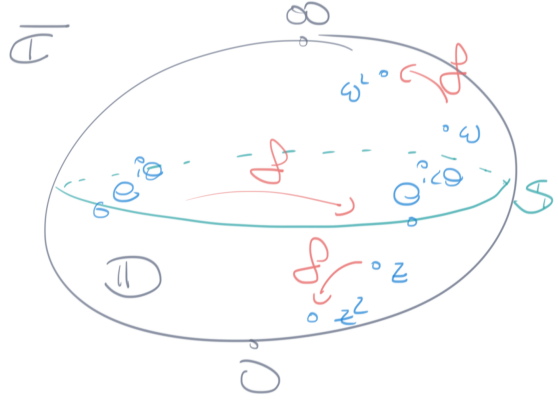


Squaring complex numbers.

Consider the transformation $f : \mathbb{C} \rightarrow \mathbb{C}$ of the complex plane defined by “squaring”, i.e.

$$f(z) = z^2.$$

The basin of attraction of the fixed point 0 is the unit disk $\mathbb{D} = \{|z| < 1\}$. Indeed, if $|z| = \lambda < 1$, then $|f^n(z)| = \lambda^{2^n} \rightarrow 0$ as $n \rightarrow \infty$. We may also extend f to an endomorphism of the Riemann sphere $\overline{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$, and then, by the same reasoning, we see that the basin of attraction of ∞ is the exterior of the disk, the set $\mathbb{D}^- = \{|z| > 1\}$. Meanwhile, it is not obvious to describe the basin of attraction of the fixed point $p = 1$. It is clear that its basin belongs to the unit circle $\mathbb{S} = \{|z| = 1\}$, and that it contains ± 1 , its square roots $\pm i$, the square roots of these points, and so on ... a countable and dense subset of the unit circle. The restriction of f to the unit circle sends $e^{2\pi i x} \mapsto e^{2\pi i 2x}$. If we identify the unit circle \mathbb{S} with \mathbb{R}/\mathbb{Z} , by means of the exponential map, we see that this restriction is the doubling map $x + \mathbb{Z} \mapsto 2x + \mathbb{Z}$. We will have much to say about it in the following.



3.4 Observables

Observables. *Observables* are functions $\varphi : X \rightarrow \mathbb{R}$ or \mathbb{C} . If the system is initially in the state x , and therefore is observed the value $\varphi(x)$ of the observable φ , after a time n observation of φ will give the value $\varphi(f^n(x))$.

Invariant functions. Particularly interesting are observables which do not change with time, that physicists call *first integrals*. The function/observable $\varphi : X \rightarrow \mathbb{R}$ is *invariant* if

$$\varphi \circ f = \varphi$$

i.e. if it is constant in each orbit. Observe that if φ is invariant, $I \subset \mathbb{R}$ and $A = \varphi^{-1}(I)$, then $f^{-1}(A) = A$. The existence of an invariant function contains the following information: if we know that $\varphi(x) = a$, then future and past of x belong to the level set $\Sigma_a = \{x \in X \text{ t.q. } \varphi(x) = a\}$, i.e. $O_f(x) \subset \Sigma_a$. Invariant functions, therefore, reduce the allowed phase space of trajectories.

Lyapunov functions. Also useful are monotone observable, which increase or decrease along trajectories, known in physics as *Lyapunov functions*. For example, if we know that $\varphi \circ f \leq \varphi$, and $\varphi(x) = a$, then the future of x does not leave the sub-level set $\Sigma_{\leq a} = \{x \in X \text{ s.t. } \varphi(x) \leq a\}$, and the past of x comes from $\Sigma_{\geq a} = \{x \in X \text{ s.t. } \varphi(x) \geq a\}$.

Energy. The energy $E(q, p) = p^2/2 + q^2/2$, which is a constant of the motion for the harmonic oscillator $\ddot{q} = -q$ (here $p = \dot{q}$), is a Lyapunov function for the damped oscillator $\ddot{q} = -\alpha\dot{q} - q$, since its time derivative is $\frac{d}{dt}E = -\alpha p^2 \leq 0$.

ex: Show that, if $\varphi : X \rightarrow \mathbb{R}$ is invariant, $I \subset \mathbb{R}$ and $A = \varphi^{-1}(I)$, then $f^{-1}(A) = A$.

ex: Show that the characteristic function of a set $A \subset X$ is invariant iff $f^{-1}(A) = A$.

Time means. The *time mean* (or *Birkhoff mean*) of the observable φ up to time $n \geq 0$ is the observable $\bar{\varphi}_n$ defined by

$$\bar{\varphi}_n(x) := \frac{1}{n+1} \sum_{k=0}^n \varphi(f^k(x))$$

i.e. the value of $\bar{\varphi}_n$ at the point x is the arithmetic mean of the values of φ along the “ n -orbit of x ”, the set $\{x, f(x), f^2(x), \dots, f^n(x)\}$. If the limit

$$\bar{\varphi}(x) = \lim_{n \rightarrow \infty} \bar{\varphi}_n(x)$$

exists, it has the meaning of “asymptotic mean value” of φ along the orbit of x . Also observe that $\bar{\varphi}(x) = (\bar{\varphi} \circ f)(x)$ at points where the limit exists.

If, in particular, 1_A denotes the characteristic function of a subset $A \subset X$, then the limit

$$\overline{1_A}(x) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \text{card} \{0 \leq k \leq n \text{ s.t. } f^k(x) \in A\}$$

if it exists, represents the “asymptotic fraction of time that the trajectory of x spend inside A ”, i.e. the asymptotic “frequency” with which the trajectory of x visit the subset A .

3.5 Invariant sets

Invariant sets. The characteristic function of a subset $A \subset X$ is invariant iff $f^{-1}(A) = A$. This motivates the following definition: a subset $A \subset X$ is *invariant* if

$$f^{-1}(A) = A$$

This condition implies that $f(A) \subset A$, and therefore a point inside an invariant set has all its history, past and future, inside the invariant set.

Observe that $O_f(x)$ is the smaller invariant set which contains x , and therefore a subset is invariant iff it is a union of big orbits. If f is invertible, $O_f(x)$ is the smaller invariant set which contains x , so that a subset is invariant iff it is a union of complete orbits, i.e. if $A = \bigcup_{x \in A} O_f(x)$.

We also say that a subset $A \subset X$ is *+invariant* (positively invariant) if $f(A) \subset A$, and *-invariant* (negatively invariant) if $f^{-1}(A) \subset A$. In particular, if A is +invariant, it is possible to define the restriction of f to A , i.e. dynamical system $f|_A : A \rightarrow A$.

ex: Discover the possible implications between the conditions

$$\begin{aligned} f^{-1}(A) = A, \quad f(A) \subset A, \quad f^{-1}(A) \subset A, \\ f(A) = A, \quad \text{e} \quad f^{-1}(A) = A = f(A) \end{aligned}$$

for a generic transformation, or a transformation which is injective, surjective, or one-to-one.

ex: Consider a set C equal to $O_f(x)$ or $O_f^+(x)$ for some $x \in X$, and determine the invariance properties of C , \overline{C} , ∂C and C' .

ex: Let $A \subset X$. Show that $\bigcup_{n \geq 0} f^n(A)$ is +invariant, indeed the smallest +invariant set which contains A .

If f is invertible, show that $\bigcup_{n \in \mathbb{Z}} f^n(A)$ is invariant, indeed the smallest invariant set which contains A .

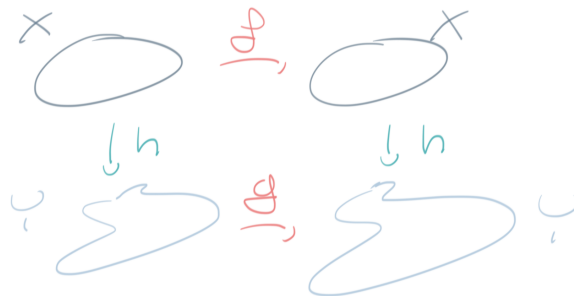
ex: Let $\varphi : X \rightarrow \mathbb{R}$ be an observable, and $A \subset X$ be the set of those points $x \in X$ such that the limit $\overline{\varphi}(x) = \lim_{n \rightarrow \infty} \overline{\varphi}_n(x)$ exists. Shows that A is invariant, and that the observable $\overline{\varphi} : A \rightarrow \mathbb{R}$ is also invariant w.r.t. the restriction $f|_A : A \rightarrow A$.

3.6 Conjugations

Conjugations. The topological dynamical systems $f : X \rightarrow X$ and $g : Y \rightarrow Y$ are (*topologically*) *conjugated* if there exists a homeomorphism $h : X \rightarrow Y$, called *conjugation*, such that

$$h \circ f = g \circ h$$

This means that arrows in the following diagram commute:



This condition may be also written $f = h^{-1} \circ g \circ h$, and is clearly an equivalence relation. By induction, we see that $f^n = h^{-1} \circ g^n \circ h$ for all times $n \geq 0$. In particular, a conjugation sends orbits of f into orbits of g , and vice-versa. The idea is that two conjugated transformations are indistinguishable from the topological point of view (we are just changing the names of the points).

e.g. Linear conjugations. Let $f : x \mapsto Ax$ be the linear map of R^N defined by the square matrix A . An automorphism $h : x \mapsto y = Ux$, defined by the invertible matrix U , defines a linear conjugation between f and the linear map $g : y \mapsto UAU^{-1}y$.

Powers and multiplications. The unit circle $\mathbb{S} \subset \overline{\mathbb{C}}$ is an invariant set for the square map $f(z) = z^2$, defined in the Riemann sphere. The exponential map $\varphi(x + \mathbb{Z}) = e^{2\pi i x}$ is a homeomorphism between the one-dimensional torus $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ and the unit circle, and defines a conjugation between the restriction $f|_{\mathbb{S}} : \mathbb{S} \rightarrow \mathbb{S}$ and the *doubling map* $E_2 : \mathbb{T} \rightarrow \mathbb{T}$, defined as $E_2(x + \mathbb{Z}) := 2x + \mathbb{Z}$.

Spirals. Consider the map $z \mapsto \lambda z$ of the Riemann sphere, where $\lambda = \rho e^{i\varphi}$ is a complex number with modulus $|\lambda| = \rho < 1$. The orbits of all points different from ∞ converge to the origin along logarithmic spirals (if the phase φ is not a multiple of 2π). As a map of the two-sphere, it is a variation of the North-South map sending $z \mapsto \rho z$. Indeed, the two are conjugated by the rotation $z \mapsto e^{i\varphi} z$.

Semi-conjugations. A continuous and onto function $h : X \rightarrow Y$ is a *semi-conjugation* between the dynamical systems $f : X \rightarrow X$ and $g : Y \rightarrow Y$ if $h \circ f = g \circ h$. In this case, g is called *factor* of f . The h -image of an orbit of f is an orbit of g , but each orbit of g may have more than one pre-image. Informally, the dynamics of f is richer than the dynamics of g . Meanwhile, when the set where h fails to be bijective is small, the two dynamics are still nearby.

4 Linear systems

The simplest higher-dimensional systems are described by linear differential equations. They provide models for the local behaviour of more general systems.

4.1 Exponential of a linear operator

Linearity & exponentials. The exponential $x(t) = e^{\lambda t}$ is the unique solution of the differential equation $\dot{x} = \lambda x$ with initial condition $x(0) = 1$. Moreover, it satisfies the functional equation $x(t+s) = x(t)x(s)$, which says that $\exp : \mathbb{R} \rightarrow \mathbb{R}^\times$ defines a homomorphism from the additive group \mathbb{R} into the multiplicative group \mathbb{C}^\times . If we try to solve a system of linear homogeneous differential equations like

$$\dot{x} = Ax,$$

with $x \in \mathbb{R}^N$ and $A \in \text{Mat}_{n \times n}(\mathbb{R})$, we are tempted to look for a solution as

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

In the following, we recall how to give a meaning to such an expression, and prove that it solves the problem. The functional equation will say that e^{tA} is a one-parameter subgroup of the general linear group $\text{GL}_N(\mathbb{R})$. The practical computation of the exponential of a matrix will make use of diagonalization, commutativity, and related considerations. More important, some qualitative aspects of solutions will derive simply from considerations on the spectrum of A , the eigenvalues of its complexification.

Exponential of a linear operator. The *exponential* of the square matrix $A = (a_{ij}) \in \text{Mat}_{N \times N}(\mathbb{C})$ is the square matrix e^A , or $\exp(A)$, defined by the power series

$$\begin{aligned} e^A &:= \sum_{k=0}^{\infty} \frac{1}{k!} A^k \\ &= I + A + \frac{1}{2}A^2 + \frac{1}{6}A^3 + \dots \end{aligned} \tag{4.1}$$

This definition makes sense because each entry of r.h.s. above is the sum of an absolutely convergent series. To see this, observe that the operator norm $\|A\| := \sup_{v \in \mathbb{C}^N, \|v\|=1} \|Av\|$ is multiplicative, i.e. satisfies $\|AB\| \leq \|A\| \|B\|$. This implies the bound

$$\|A^k/k!\| \leq \|A\|^k/k!$$

There follows, since all norms in a finite dimensional vector space are equivalent, that the absolute value of each entry of the series (4.1) is bounded by a constant times the convergent series $\sum_{k=0}^{\infty} \|A\|^k/k! = e^{\|A\|}$. Bytheway, this also implies the bound

$$\|e^A\| \leq e^{\|A\|}.$$

It is clear that if the matrix A is real, then also its exponential e^A is real.

If A and B are similar matrices, i.e. $A = U^{-1}BU$ for some $U \in \text{GL}_N(\mathbb{C})$, then also their exponentials are similar, since powers of A are $A^n = U^{-1}B^nU$ for all $n \geq 0$, and therefore one easily justifies the following computation

$$\begin{aligned} e^A &= I + U^{-1}BU + \frac{1}{2}U^{-1}B^2U + \dots \\ &= U^{-1} \left(I + B + \frac{1}{2}B^2 + \dots \right) U \\ &= U^{-1}e^B U. \end{aligned} \tag{4.2}$$

Therefore, if L is a linear operator defined in a finite-dimensional vector space isomorphic to \mathbb{C}^N or \mathbb{R}^N , represented in some fixed basis by the matrix A , then the formula (4.1) defines a linear operator

$$e^L = I + L + \frac{1}{2}L^2 + \frac{1}{6}L^3 + \dots$$

According to formula (4.2), this definition does not depend on the chosen basis.

Exponential of diagonalizable matrices. If A is a diagonal matrix with eigenvalues λ_k 's, i.e.

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N) := \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_N \end{pmatrix}$$

(missing entries are zero) then a straightforward computation shows that its exponential is also diagonal, and indeed

$$e^\Lambda = \text{diag}(e^{\lambda_1}, \dots, e^{\lambda_N}) = \begin{pmatrix} e^{\lambda_1} & & & \\ & e^{\lambda_2} & & \\ & & \ddots & \\ & & & e^{\lambda_N} \end{pmatrix}.$$

In particular, if A is diagonalizable, i.e. $A = U^{-1}\Lambda U$ with Λ diagonal and $U \in \text{GL}_N(\mathbb{C})$, then its exponential is similar to the diagonal matrix e^Λ , namely $e^A = U^{-1}e^\Lambda U$. Thus, exponentials of diagonalizable matrices are easy to compute, provided we know the change of coordinates U that diagonalizes the matrix.

An important consequence is a relation between the exponential and the principal invariants of a square matrix, the determinant and the trace. It says that

$$\boxed{\det(e^A) = e^{\text{tr}A}} \quad (4.3)$$

This formula is obvious if A is diagonalizable, and follows by continuity in the general case, because the set of diagonalizable matrices is dense in the space $\text{Mat}_{n \times n}(\mathbb{C})$ of complex square matrices (a generic degree n complex polynomial has n distinct roots).

ex: Show that if v is an eigenvector of the linear operator L with eigenvalue λ , then v is also an eigenvector of e^L , with eigenvalue e^λ .

One-parameter groups of matrices. Given a matrix $A \in \text{Mat}_{N \times N}(\mathbb{C})$, we may consider the family of matrices

$$G(t) := e^{tA},$$

parametrized by a “time” $t \in \mathbb{R}$. It is clear that $G(0) = I$. The series of functions $t \mapsto (e^{tA})_{ij}$ which define the entries of e^{tA} converge uniformly in any bounded interval of the real line, as well the series of their derivatives. In particular, the time derivatives may be computed term-wise. The result is that

$$\frac{d}{dt} G(t) = \sum_{k=0}^{\infty} \frac{t^k}{k!} A^{k+1} = A G(t) = G(t) A \quad (4.4)$$

In particular, A commutes with $G(t)$.

The derivative of $F(t) := e^{tA}e^{-tA}$ is equal, by the Leibniz rule applied to every entry of the product, to $F'(t) = AF(t) - F(t)A = 0$, because A commutes with $G(t)$. By the mean value theorem, $F(t) = F(0) = I$. Therefore, $G(t) = e^{tA}$ is invertible, and its inverse is $(e^{tA})^{-1} = e^{-tA}$. Thus, the exponential sends $\exp : \text{Mat}_{n \times n}(\mathbb{C}) \rightarrow \text{GL}_n(\mathbb{C})$.

Theorem 4.1. Let $A \in \text{Mat}_{N \times N}(\mathbb{C})$. The unique solution of the linear differential equation

$$\dot{X} = AX \quad \text{or} \quad \dot{X} = XA,$$

with initial condition $X(0) = X_0 \in \text{GL}(N, \mathbb{C})$, is

$$X(t) = e^{tA} X_0 \quad \text{or} \quad X(t) = X_0 e^{tA},$$

respectively.

Proof. It is clear, by the above computation, that $e^{tA}X_0$ or X_0e^{tA} are solutions of the two problems. In order to prove uniqueness, we may observe that if $X(t)$ is a solution, then the matrix $X(t)e^{-tA}$ (or $e^{-tA}X(t)$ in the second case) does not depend on time, since its derivative is zero, and therefore is constant and equal to its initial value X_0 . \square

Observe that the two differential equations in the above theorem are not the same, the product between matrices does not commute, in general. Indeed, if A and B do not commute, the three exponentials e^{A+B} and e^Ae^B and e^Be^A may all be different from each other. What is true is the following.

Theorem 4.2. *If A and B commute, i.e. if $AB = BA$, then*

$$e^{A+B} = e^Ae^B = e^Be^A.$$

Proof. If A commutes with B , then all its powers A^k also commute with all the powers B^j , and therefore with the exponentials e^{tB} and e^{tA} , and viceversa. There follows that the derivative of

$$H(t) = e^{t(A+B)} - e^{tA}e^{tB}$$

is, using formulas 4.4,

$$H'(t) = (A+B)e^{t(A+B)} - Ae^{tA}e^{tB} - e^{tA}e^{tB}B = (A+B)H(t)$$

By the uniqueness theorem 4.1, $H(t) = e^{t(A+B)}H(0)$. But $H(0) = 0$, therefore $H(t) = 0$ for all times t , and in particular for $t = 1$. \square

In particular, since all multiples tA of A commute, the family of the $G(t) = e^{tA}$, with $t \in \mathbb{R}$, is a *one-parameter subgroup* of the general linear group $GL_N(\mathbb{C})$, i.e. satisfies

$$e^{0A} = I \quad \text{and} \quad e^{tA}e^{sA} = e^{(t+s)A}.$$

In other words, the correspondence $t \mapsto e^{tA}$ is an homomorphism of the additive group \mathbb{R} into $GL_N(\mathbb{C})$. Its image is a curve in the linear group, which passes through the identity for $t = 0$, and solves the differential equation $\dot{G} = AG$. The matrix A is called *generator* of the subgroup $\{G(t)\}_{t \in \mathbb{R}}$, and may be obtained as the derivative

$$A = \dot{G}(0) = \lim_{t \rightarrow 0} \frac{G(t) - I}{t}.$$

Thus, A is the velocity of the curve $G(t)$ at time $t = 0$.

4.2 Linear flows

Linear systems. A *homogeneous linear system* with constant coefficients is an autonomous differential equation

$$\dot{x} = L(x) \tag{4.5}$$

for $x(t) \in \mathbb{R}^N$, defined by a linear vector field $L \in \text{End}(\mathbb{R}^N)$. The origin is an equilibrium solution, since $L(0) = 0$ by linearity. Fixed a basis of \mathbb{R}^N , e.g. the canonical basis, the system may be written in matrix notation as

$$\dot{x} = Ax,$$

where $x(t) = (x_1(t), x_2(t), \dots, x_n(t))^T \in \mathbb{R}^N$ is a column vector, $A = (a_{ij}) \in \text{Mat}_{N \times N}(\mathbb{R})$ is the matrix which represents the linear vector field L in the chosen basis, and Ax denotes the usual product between matrices. By the proof of theorem 4.1, the solution with initial condition $x(0) = x_0 \in \mathbb{R}^N$ is given by

$$x(t) = e^{tA}x_0.$$

The flow of the linear vector field L is the one-parameter group of linear maps $\Phi_t = e^{tL}$, given, in the chosen basis, by $\Phi_t(x) := e^{tA}x$.

Thus, if we want to understand solutions of a linear system, we must compute the exponential of the linear vector field, in some convenient basis.

Diagonalizable linear systems. Assume that A is diagonalizable, and has n real eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_N$ (not necessarily distinct) with linearly independent eigenvectors v_1, v_2, \dots, v_N , respectively, so that $Av_k = \lambda_k v_k$ and the v_k 's form a basis of \mathbb{R}^N . Then the solution of (4.5) with initial conditions $x(0) = \sum_k a_k v_k$ is a superposition

$$x(t) = \sum_{k=1}^N e^{t\lambda_k} a_k v_k$$

The qualitative asymptotic behaviour of solutions is therefore decided by the signs of the eigenvalues.

For example, if all the eigenvalues are negative, i.e. $\lambda_k < 0$ for all $k = 1, \dots, N$, then all solutions decay, exponentially fast to the origin, i.e.

$$\|x(t)\| \leq e^{-\alpha t} \|x(0)\|$$

for some $\alpha = \min_k |\lambda_k| > 0$. The origin is then an “asymptotically stable” equilibrium, or a “sink”.

If, on the other side, all the eigenvalues are positive, i.e. $\lambda_k > 0$ for all $k = 1, \dots, N$, then all solutions different from the equilibrium solution diverge exponentially fast, i.e.

$$\|x(t)\| \geq e^{\beta t} \|x(0)\|$$

for some $\beta = \min_k \lambda_k > 0$. The origin is an “asymptotically unstable equilibrium”, or a “source”.

More interesting is the mixed situation of a saddle, with some stable directions and some unstable directions. The case with some zero eigenvalue, i.e. some indifferent directions, is clearly non generic, although physically interesting (the harmonic oscillator is such a case!).

On the other side, generic real matrices are not diagonalizable. To understand their exponentials, we must complexify and use the Jordan normal form.

Complexification. The *complexification* of the real vector space \mathbb{R}^N is the complex vector space $\mathbb{C}^N := \mathbb{R} \oplus i\mathbb{R}$, i.e. the set of vectors $z = x \oplus iy \approx x + iy$, with $x, y \in \mathbb{R}^N$, equipped with the natural sum and multiplication by complex scalars.

The complexification of the linear map $x \mapsto L(x)$ defined, in the canonical basis of \mathbb{R}^N , by a matrix $A \in \text{Mat}_{n \times n}(\mathbb{R})$ according to $x \mapsto Ax$, is the linear operator $z \mapsto L^{\mathbb{C}}(z)$ defined by the same matrix, i.e. according to $z = x + iy \mapsto Az = Ax + iAy$.

The *spectrum* of the linear operator L (in a finite dimensional linear space) is the set $\sigma(L) \subset \mathbb{C}$ of the eigenvalues of its complexification $L^{\mathbb{C}}$, i.e. complex roots of the characteristic polynomial $P_A(t) := \det(t - A)$. By Gauss' fundamental theorem of arithmetic, the characteristic polynomial factorizes as a product

$$P_A(t) = \prod_{\lambda \in \sigma(L)} (t - \lambda)^{m_\lambda}.$$

The integer exponent m_λ is called (*algebraic*) *multiplicity* of the eigenvalue λ . It is clear that $\sum_{\lambda \in \sigma(A)} m_\lambda = n$.

Complexification in dimension two. The relevant example, for our purposes, is the following. Let $x \mapsto L(x)$ be the linear operator defined, in the canonical basis $e_1 = (1, 0)$ and $e_2 = (0, 1)$ of \mathbb{R}^2 , by the real matrix

$$A = \begin{pmatrix} \alpha & \omega \\ -\omega & \alpha \end{pmatrix}$$

Thus, $L(e_1) = \alpha e_1 - \omega e_2$ and $L(e_2) = \omega e_1 + \alpha e_2$. Then the complexified operator $z \mapsto L^{\mathbb{C}}(z)$ is defined, in the basis $v_+ = e_1 + ie_2$ and $v_- = e_1 - ie_2$ of \mathbb{C}^2 , by the diagonal matrix

$$\Lambda = \begin{pmatrix} \lambda & 0 \\ 0 & \bar{\lambda} \end{pmatrix}$$

where $\lambda = \alpha + i\omega$. Indeed, a computation shows that $A(e_1 + ie_2) = (\alpha + i\omega)(e_1 + ie_2)$ and $A(e_1 - ie_2) = (\alpha - i\omega)(e_1 - ie_2)$.

Vice-versa, let $L^{\mathbb{C}}$ be the complexification of a real operator defined by the real two-by-two matrix A , and let v_+ be an eigenvector of $L^{\mathbb{C}}$ with eigenvalue $\lambda = \alpha + i\omega$, so that, in the canonical basis, $Av_+ = \lambda v_+$. Then $v_- := \overline{v_+}$ is an eigenvector of $L^{\mathbb{C}}$ with eigenvalue $\bar{\lambda} = \alpha - i\omega$. Indeed, since the entries of A are real, the roots of the characteristic polynomial comes in pairs of conjugated complex numbers, and one check that $Av_- = A\overline{v_+} = \overline{Av_+} = \overline{\lambda v_+} = \bar{\lambda} v_-$. There follows that, in the real basis $e_+ = (v_+ + v_-)/2$ and $e_- = (v_+ - v_-)2i$, which is therefore a basis of the real vector space $\mathbb{R}^2 \subset \mathbb{C}^2$, the real operator L is represented by the matrix A as above.

So, a diagonalizable complexified real linear operator in the plane with a couple of complex conjugate eigenvalues $\lambda_{\pm} = \alpha \pm i\omega$ “corresponds” to a two-by-two real matrix which is the sum of a multiple of the identity αI and an anti-symmetric matrix Ω as below

$$\alpha I + \Omega := \alpha \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + \begin{pmatrix} 0 & \omega \\ -\omega & 0 \end{pmatrix}.$$

Since any matrix commute with any multiples of the identity, by theorem 4.2 we may compute separately the exponentials of $t\rho I$ and $t\Omega$, and then multiply the results. The flow of the diagonal part is simply $e^{t\alpha I} = e^{\alpha t}I$. A computation (using the power series of the trigonometric functions $\sin t$ and $\cos(t)$) shows that the flow defined by the antisymmetric matrix Ω above is

$$e^{t\Omega} = \begin{pmatrix} \cos(\omega t) & \sin(\omega t) \\ -\sin(\omega t) & \cos(\omega t) \end{pmatrix}$$

i.e. it is a clockwise rotation $R_{t\omega}$ by an angle $t\omega$. Multiplying, we finally get

$$e^{tA} = e^{t(\alpha I + \Omega)} = e^{\alpha t} \begin{pmatrix} \cos(\omega t) & \sin(\omega t) \\ -\sin(\omega t) & \cos(\omega t) \end{pmatrix}$$

So, the flow of A is a rotation with angular frequency ω (or frequency $\nu = \omega/2\pi$) together with stretching/contraction with exponential rate α . Orbits are logarithmic spirals entering or coming from the origin, depending on the sign of α .

The case $\alpha = 0$ corresponds to pure rotations (this is the case of the harmonic oscillator $\ddot{x} = -\omega^2 x$).

4.3 Linear systems in the plane

Linear systems in the plane. We have now all the tools to understand the general linear system of differential equations

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

in the plane \mathbb{R}^2 , defined by a real 2×2 matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}.$$

Let λ_+ and λ_- be the eigenvalues of the complexification of A , i.e. the complex roots (possibly equal) of the characteristic polynomial $\det(tI - A)$. The product $\lambda_+ \lambda_-$ of the eigenvalues is $q = \det(A) = ad - bc$, and the sum $\lambda_+ + \lambda_-$ of the eigenvalues is $p = \text{tr}(A) = a + d$. Eigenvalues are therefore

$$\lambda_{\pm} = \frac{p \pm \sqrt{\Delta}}{2},$$

where the “discriminant” is $\Delta = p^2 - 4q$.

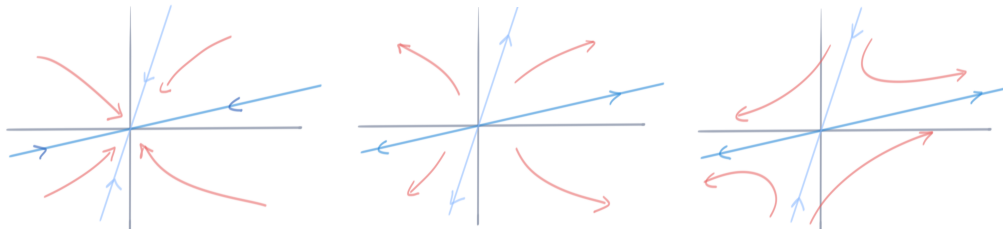
Two independent eigenvectors. If the matrix A is diagonalizable over the reals, i.e. admits two linearly independent eigenvectors with real eigenvalues $\lambda_{\pm} \in \mathbb{R}$ (possibly equal), then the system is linearly equivalent to

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} \lambda_+ & 0 \\ 0 & \lambda_- \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

Solutions are

$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = \begin{pmatrix} e^{\lambda_+ t} & 0 \\ 0 & e^{\lambda_- t} \end{pmatrix} \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$

The origin is called *stable node* if $\lambda_{\pm} < 0$, *unstable node* if $\lambda_{\pm} > 0$, or *saddle* if $\lambda_- < 0 < \lambda_+$.



Stable node, unstable node and saddle.

Only one eigenvector. Assume that the matrix A admits just one eigenvector v (or better, a one-dimensional space of eigenvectors), with eigenvalue $\lambda \in \mathbb{R}$. Then the matrix representing the operator in a basis v, w (where w is any other linearly independent vector) is upper triangular, with both diagonal entries equal to λ (for otherwise a different second diagonal entry would be another eigenvalue and would therefore yield a second independent eigenvector) and a non-zero upper right entry (for otherwise w would be a second, linearly independent eigenvector, and the operator diagonalizable). Therefore, the operator sends $v \mapsto \lambda v$ and $w \mapsto av + \lambda w$ for some $a \neq 0$. Thus, in the basis formed by v and w/a , the operator is defined by the upper triangular matrix

$$\begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix} \quad (4.6)$$

This shows that the system is linearly equivalent to

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} \lambda & 1 \\ 0 & \lambda \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

Since (4.6) is a sum of the diagonal matrix ρI and the nilpotent matrix $N = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$, which indeed satisfies $N^2 = 0$, the power series defining its exponential is actually a polynomial of first degree, and solutions are

$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = e^{\lambda t} \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$

The origin is called *degenerate node*, stable or unstable, depending on the sign of ρ .

ex: A more conceptual proof of the above observation introduces to the so called “Jordan chains”, the building blocks of the Jordan normal form of a matrix. Let L be an operator on \mathbb{R}^2 , and assume that it admits only one-dimensional eigenspace, say generated by the eigenvector v_1 with eigenvalue λ . This means that the kernel of the operator $N = L - \lambda$ is one-dimensional. For dimensional reasons, also the range of N is one-dimensional. Their intersection cannot be trivial, for otherwise the whole space would be a direct sum $\mathbb{R}^2 = \text{ran}(N) \oplus \text{kernel}(N)$ of two one-dimensional invariant subspaces, and therefore the operator L would be diagonalizable. Thus, there exists a non-zero vector v_0 such that $Nv_0 = v_1$. The vectors v_0 and v_1 are independent. Indeed, let $av_0 + bv_1 = 0$. Applying N , we get $av_1 = 0$, since $Nv_1 = 0$ and $Nv_0 = v_1$, and therefore $a = 0$. But then $bv_1 = 0$ implies that also $b = 0$ (the vectors v_0 and $v_1 = Nv_0$ form a so called “Jordan chain”). Finally, one just observes that, in the basis formed by v_1 and v_0 , the operator L is represented by the matrix (4.6), since $Lv_1 = \lambda v_1$ and $Lv_0 = v_1 + \lambda v_0$.

Complex eigenvalues. If the matrix A has no real eigenvalue, then its complexification admits two complex conjugate eigenvalues λ and $\bar{\lambda}$. If the eigenvalues are purely imaginary, say $\lambda_{\pm} = \pm i\omega$, with $\omega > 0$, then, by the previous discussion, the system is linearly equivalent to

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} 0 & \omega \\ -\omega & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

This is an harmonic oscillator $\ddot{x} = -\omega^2 x$ with angular frequency ω . Solutions are

$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = \begin{pmatrix} \cos(\omega t) & \sin(\omega t) \\ -\sin(\omega t) & \cos(\omega t) \end{pmatrix} \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$

Orbits are ellipsis, and the origin is called (*indifferent*) *focus*. Trajectories which start near the origin stay near the origin for all times, still not being asymptotic to the origin.

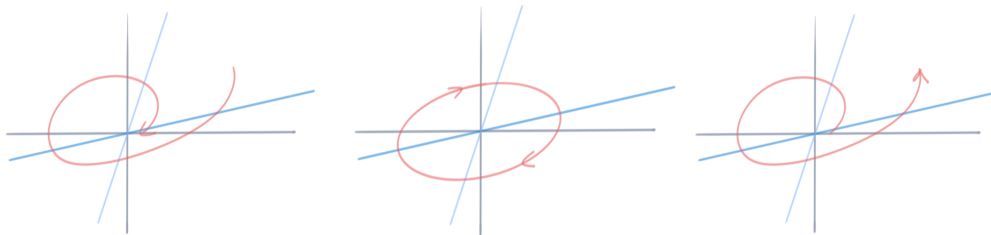
The generic case is a complexified matrix with complex eigenvalues $\lambda_{\pm} = \rho \pm i\omega$, with real part $\rho \neq 0$. The system is linearly equivalent to

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} \rho & \omega \\ -\omega & \rho \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

Solutions are

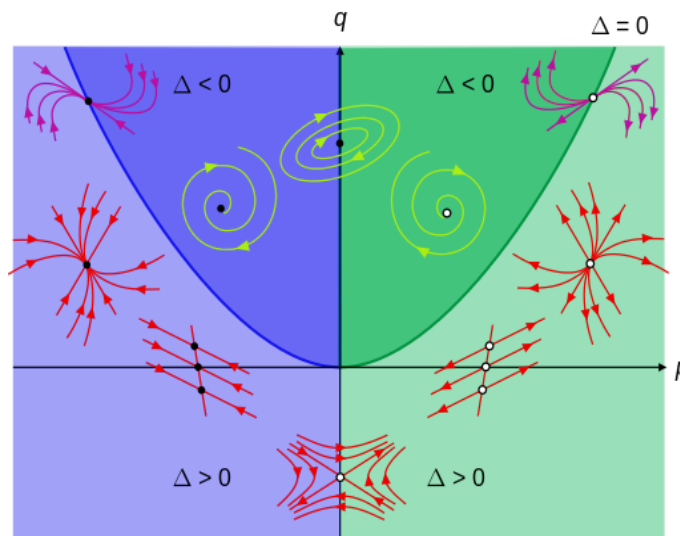
$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = e^{\rho t} \begin{pmatrix} \cos(\omega t) & \sin(\omega t) \\ -\sin(\omega t) & \cos(\omega t) \end{pmatrix} \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$

orbits are logarithmic spirals that comes or enter into the origin, depending on the sign of ρ . The origin is called *unstable focus* if $\rho > 0$, or *stable focus* if $\rho < 0$.



Stable, indifferent and unstable focus.

Global picture. It is clear that the stability or unstability of nodes or foci is preserved under small perturbations of the parameters (the entries of the matrix A). Here is a famous picture of the different phase portraits, depending on the trace and determinant of the matrix.



$$\begin{aligned} \frac{dx}{dt} &= Ax + By & p &= A + D \\ \frac{dy}{dt} &= Cx + Dy & q &= AD - BC \\ & & \Delta &= p^2 - 4q \end{aligned}$$

By Maschen, from [Wikimedia Commons](#).

ex: Discuss the degenerate cases when one of the eigenvalues is zero (so that the matrix A is not invertible).

ex: Consider the “inverted oscillator”

$$\begin{aligned}\dot{q} &= p \\ \dot{p} &= q\end{aligned}$$

Find the nature of the equilibrium, and determine the generic solution.

ex: Sketch the phase portrait (i.e. some orbits near the equilibrium in the phase space) of the following linear systems.

$$\begin{array}{ccc}\left\{ \begin{array}{l} \dot{x} = x - y \\ \dot{y} = x + y \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = 2x + y \\ \dot{y} = x + y \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = 4x \\ \dot{y} = 2x - y \end{array} \right. \\ \left\{ \begin{array}{l} \dot{x} = 6x + 5y \\ \dot{y} = x + 2y \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = -x + 2y \\ \dot{y} = 3y \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = -7x + y \\ \dot{y} = -4x - 3y \end{array} \right. \\ \left\{ \begin{array}{l} \dot{x} = y \\ \dot{y} = -4x \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = -x + 5y \\ \dot{y} = -5x - y \end{array} \right. & \left\{ \begin{array}{l} \dot{x} = x + 5y \\ \dot{y} = -5x + y \end{array} \right.\end{array}$$

ex: The current $I(t)$ in a LRC circuit is a solution of the homogeneous differential equation

$$L\ddot{I} + R\dot{I} + \frac{1}{C}I = 0$$

Write the corresponding linear system for $x(t) = I(t)$ and $y(t) = \dot{I}(t)$, and sketch the possible phase portraits, depending on the relative values of the positive parameters L , R and C .

4.4 Jordan normal form

In the higher-dimensional case, the useful normal form to understand exponentials is the Jordan normal form.

Generalized eigenspaces. Let $L : \mathbb{C}^N \rightarrow \mathbb{C}^N$ linear operator defined in a complex linear space \mathbb{C}^N . Given a scalar λ , let L_λ denotes the operator $L - \lambda$. If the kernel of L_λ is not trivial, then λ is an eigenvalue of L , and $V_\lambda = \text{kernel}(L_\lambda)$ is its associated proper space, made of eigenvectors v such that $Lv = \lambda v$.

A non-zero vector $v \in \mathbb{C}^N$ is said *generalized eigenvector* if it is in the kernel of some power of L_λ , i.e. if there exists $\lambda \in \mathbb{C}$ and some minimal integer $m \geq 1$ such that $L_\lambda^m v = 0$. The non-zero integer m is called *period* of v , and the vector v itself is also called *L_λ -cyclic* (meaning that the orbit of v by the map L_λ is formed by m distinct non-zero vector).

If the period is $p = 1$, then v is an eigenvector of L . In general, the m vectors

$$v_1 = L_\lambda^{m-1}v \quad v_2 = L_\lambda^{m-2}v \quad \dots \quad v_m = v \quad (4.7)$$

are all generalized eigenvectors, since $L_\lambda^k v_k = L_\lambda^k L_\lambda^{m-k} v = L_\lambda^m v = 0$, and the first one, v_1 , is an eigenvector of L with eigenvalue λ .

Theorem 4.3. *If v is a L_λ -cyclic vector of period m , then the m vectors (4.7) are linearly independent and generate a L -invariant subspace of generalized eigenvectors.*

Proof. If $a_1 v_1 + a_2 v_2 + \dots + a_m v_m = 0$ for some non-zero vector (a_1, a_2, \dots, a_m) , then $p(L_\lambda)v = 0$ where $p(t)$ is the non-zero polynomial

$$p(t) = a_1 t^{m-1} + a_2 t^{m-2} + \dots + a_{m-1} t + a_m$$

But also $q(L_\lambda)v = 0$, where $q(t) = t^m$, because v is L_λ -cyclic with period m . If $h(t)$ denotes the maximum common divisor between the polynomials $p(t)$ and $q(t)$, then there exist polynomials $f(t)$ and $g(t)$ such that $h(t) = f(t)p(t) + g(t)q(t)$. Therefore, also $h(L_\lambda)v = 0$. But $h(t)$ is a power of t (since it divides t^m) of degree $\deg(h) = k \leq m-1$ (since it divides $p(t)$). Thus, $h(t) = t^k$,

and therefore $L_\lambda^k v = 0$. This contradicts the fact that m is the period of v . Thus, the vectors are linearly independent.

The v_k 's are all generalized eigenvectors, because $L_\lambda^k v_k = L_\lambda^k L_\lambda^{m-k} v = L_\lambda^m v = 0$. Finally, the subspace generated by the v_k 's is L -invariant, because

$$Lv_k = L(L_\lambda^{m-k} v) = L_\lambda^{m-k+1} v + \lambda L_\lambda^{m-k} v = v_{k-1} + \lambda v_k$$

where, clearly, we set $v_0 = (L - \lambda I)^n v = 0$. \square

The kernels $\ker(L_\lambda^k)$ are called *generalized eigenspaces* of order k , and one easily sees that $\ker(L_\lambda) \subset \ker(L_\lambda^2) \subset \cdots \subset \ker(L_\lambda^n)$. Moreover, if λ is an eigenvalue of L , then the space of generalized eigenvectors associated to the eigenvalue λ is equal to $\ker(L_\lambda^n)$ (which, of course, may coincide with $\ker(L_\lambda^k)$ for some smaller $k \leq n$).

Jordan blocks. If it happens that the vectors (4.7) span the whole space \mathbb{C}^N , i.e. if $m = N$, then the entire space is said *cyclic*. The computation in the proof above shows that

$$Lv_1 = \lambda v_1 \quad \text{and} \quad Lv_k = \lambda v_k + v_{k-1} \quad \text{for } 2 \leq k \leq N.$$

Therefore, the matrix which represents the linear operator L in this basis v_1, v_2, \dots, v_N is

$$J_\lambda = \begin{pmatrix} \lambda & 1 & & & \\ & \lambda & 1 & & \\ & & \ddots & \ddots & \\ & & & \lambda & 1 \\ & & & & \lambda \end{pmatrix} \quad (4.8)$$

In particular, v_1 is the unique eigenvector, with eigenvalue λ , the proper space $V_\lambda = \ker(L_\lambda)$ being the line $\mathbb{C}v_1$. Thus, the geometric multiplicity of λ is equal to 1. The matrix (4.8) is called *Jordan block* of dimension N , and the basis (4.7) is called *Jordan basis*, or *Jordan chain* of length N . The vector v_N is called *generator*, or *lead vector* of the Jordan chain.

Observe that a Jordan block of dimension n has the form

$$J_\lambda = \lambda I + D$$

where D is the nilpotent (upper triangular) matrix

$$D = \begin{pmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \ddots & \ddots & \\ & & & 0 & 1 \\ & & & & 0 \end{pmatrix} \quad (4.9)$$

which satisfies $De_k = e_{k-1}$, if the e_k 's denote the column vectors of the canonical basis of \mathbb{C}^N , and $D^N = 0$.

The characteristic polynomial of a Jordan block J of length N is $P_J(z) = (z - \lambda)^N$, and therefore the algebraic multiplicity of the eigenvalue λ is N . The minimal polynomial (the monic polynomial of minimal degree such that $f(J) = 0$) is also $M_J(z) = (z - \lambda)^N$ (to be compared with the minimal polynomial of the diagonalizable matrix $\Lambda = \lambda I$ of order N , which is only $M_\Lambda(z) = (z - \lambda)$).

e.g. Derivative and quasi-polynomials. The paradigmatic example is the derivative operator ∂ , defined by $(\partial f)(t) := f'(t)$ on complex-valued functions $f(t)$ of a real variable t . Its eigenfunctions are the exponentials $e^{\lambda t}$, since $\partial(e^{\lambda t}) = \lambda e^{\lambda t}$. On the other side, it is nilpotent on the space of polynomials $p(t)$ of fixed degree, say $\deg(p) < n$, where $\partial^n = 0$. Exponentials and polynomials combine to form the spaces $Q_{\lambda,n} \approx \mathbb{C}^N$ of quasi-polynomials $f(t) = p(t)e^{\lambda t}$, where λ is a fixed complex exponent and the $p(t)$'s are polynomials of $\deg(p) < n$. These are cyclic spaces for the derivative, since $(\partial - \lambda)^n = 0$, and a generating vector is $t^{n-1}e^{\lambda t}$. The eigenvector of ∂ is, of course, $f(t) = e^{\lambda t}$, and has eigenvalue λ . A Jordan basis is

$$e^{\lambda t} \quad t e^{\lambda t} \quad \frac{1}{2} t^2 e^{\lambda t} \quad \dots \quad \frac{1}{(n-1)!} t^{n-1} e^{\lambda t}$$

In this basis, the operator ∂ is represented by the matrix (4.8).

Flow of a Jordan block. The exponential of t times a Jordan block J_λ of dimension N , which defines the flow of the linear differential equation

$$\dot{v} = J_\lambda v$$

defined in a cyclic space, is easily computed. Indeed, since D commute with λI , we may compute separately the two exponentials and then multiply. But since D is nilpotent, namely $D^N = 0$, the series defining the exponential terminates, and indeed

$$e^{tD} = I + tD + \frac{t^2}{2}D^2 + \cdots + \frac{t^{N-1}}{(N-1)!}D^{N-1}.$$

So, the exponential of tJ_λ is simply

$$e^{tJ_\lambda} = e^{t\lambda} \left(I + tD + \frac{t^2}{2}D^2 + \cdots + \frac{t^{N-1}}{(N-1)!}D^{N-1} \right).$$

It is clear, since polynomials corrections are negligible compared with exponential growth or decay, that the asymptotic behaviour of solutions of the linear system $\dot{v} = J_\lambda v$ only depends on the sign of the real part of λ , provided it is not zero.

Theorem 4.4. *If $\Re(\lambda) < 0$, then for all $0 < \alpha < |\Re(\lambda)|$ there exists a constant C such that*

$$\|e^{tJ_\lambda} v\| \leq C e^{-\alpha t} \|v\| \quad \text{for } t \geq 0.$$

Proof. If we write a generic vector as a superposition $v = \sum_k a_k v_k$ of the vectors v_k 's of the Jordan basis, we see that

$$e^{tJ_\lambda} v = e^{\lambda t} \left(\sum_i a_i p_{ik}(t) \right) v_k$$

where the $p_{ik}(t)$'s are certain polynomials of degree $< n$, which only depend on the dimension n of the Jordan block. Assume that $\Re(\lambda) = -\rho < 0$. Take any $0 < \alpha < \rho$, and set $\varepsilon = \rho - \alpha > 0$. We may define a norm on the cyclic space according to $\|v\|_\lambda := \max_k |a_k|$. Then, if $M = \max_{i,k} M_{ik}$ denotes the maximal value of the $M_{ik} = \sup_{t \geq 0} |e^{-t\varepsilon} p_{ik}(t)|$, we clearly have

$$\|e^{tJ_\lambda} v\|_\lambda \leq M e^{-\alpha t} \|v\|_\lambda$$

for all $t \geq 0$. Since all norms in a finite dimensional vector space are equivalent, this finally implies that claimed inequality, for some other constant C , holds for the standard or any other norm in the cyclic space. \square

Thus, if $\Re(\lambda) < 0$, all vectors are exponentially contracted by the flow of J_λ , and decay to zero exponentially fast as $t \rightarrow \infty$.

Reversing the arrow of time, one shows that if $\Re(\lambda) = \rho > 0$ and $\rho > \beta > 0$, then there exists a constant C such that

$$\|e^{-tJ_\lambda} v\| \leq C e^{-\beta t} \|v\| \quad \forall t \geq 0$$

Thus, if $\Re(\lambda) > 0$, all vectors are exponentially stretched by the flow of J_λ , and decay to zero exponentially fast as $t \rightarrow -\infty$.

Jordan normal form. It happens that any linear operator in a finite dimensional complex vector space is a direct sum of Jordan blocks.

Theorem 4.5 (Jordan normal form). *Let L be a linear operator in a finite-dimensional complex vector space \mathbb{C}^N . The total space splits as a direct sum $\mathbb{C}^N = E_{\lambda_1} \oplus E_{\lambda_2} \oplus \cdots \oplus E_{\lambda_d}$ of cyclic L -invariant subspaces.*

Therefore, if we chose a Jordan basis in any invariant cyclic subspace E_{λ_k} , the matrix that represents the linear operator L in the resulting basis is block diagonal as

$$J = \begin{pmatrix} J_{\lambda_1} & 0 & \dots & 0 \\ 0 & J_{\lambda_2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & J_{\lambda_d} \end{pmatrix} \quad (4.10)$$

where each $J_{\lambda_k} = \lambda_k I + D_k$ is a Jordan block as (4.8). The λ_k 's are the eigenvalues of L , the roots of the characteristic polynomial $P_A(z) = \det(zI - A)$, where A is the matrix that represents L in the canonical basis. Indeed, the characteristic polynomial factorizes as a product $P_A(z) = \prod_{\lambda \in \sigma(A)} (z - \lambda)^{m_\lambda}$, where m_λ is the (*algebraic multiplicity*) of the eigenvalue λ , which is equal to the sum of the dimensions of the Jordan blocks with $\lambda_k = \lambda$, i.e. to the dimension of the generalized eigenspace $\ker(L_\lambda^n)$. The (*geometric multiplicity*) of the eigenvalue λ is the dimension of the proper space $\ker(L_\lambda)$, which is equal to the cardinality of those Jordan blocks with $\lambda_k = \lambda$. The minimal polynomial of A is a product $M_A(z) = \prod_{\lambda \in \sigma(A)} (z - \lambda)^{\mu_\lambda}$, where μ_λ is the dimension of the largest Jordan block with $\lambda_k = \lambda$.

If A is the matrix that represents the linear operator L in the canonical basis (or in any other basis), then there exists an invertible matrix $G \in \text{GL}_N(\mathbb{C})$ (whose columns are the vectors of the Jordan bases) such that $G^{-1}AG = J$. The *canonical form* J is unique modulo permutations of the blocks. In particular, the matrix A may be represented as a sum

$$A = \Lambda + D$$

of a *semi-simple*, i.e. diagonalizable, matrix $\Lambda = G(\lambda_1 I \oplus \lambda_2 I \oplus \dots)G^{-1}$ and a nilpotent matrix $D = G(D_1 \oplus D_2 \oplus \dots)G^{-1}$ which commute, i.e. such that $\Lambda D = D\Lambda$.

Clear proofs of the Jordan normal form theorem 4.5 can be found in the classical [HS74], or in any good reference on linear algebra, as for example [La87, Ax97].

Normal form of real operators. We now consider a linear operator $L : \mathbb{R}^N \rightarrow \mathbb{R}^N$, defined, in the canonical basis, by a matrix $A \in \text{Mat}_{n \times n}(\mathbb{R})$. We may think at A as a complex matrix, representing the complexified operator $L^\mathbb{C} : \mathbb{C}^N \rightarrow \mathbb{C}^N$, and as such conjugated to a block diagonal matrix as (4.10) above. Eigenvalues are real, or come in couples of complex conjugated pairs $\lambda_\pm = \alpha \pm i\omega$, since the characteristic polynomial has real coefficients.

Theorem 4.6 (Jordan normal form for real operators). *Let L be a linear operator on the real vector space \mathbb{R}^N . The total space splits as a direct sum of invariant subspaces E_λ or $E_{\lambda, \bar{\lambda}}$, namely*

$$\mathbb{R}^N = \left(\bigoplus_{\lambda \in \mathbb{R}} E_\lambda \right) \oplus \left(\bigoplus_{\lambda \in \mathbb{C} \setminus \mathbb{R}} E_{\lambda, \bar{\lambda}} \right),$$

where the operator is represented by a matrix of the form (4.8), for some real eigenvalue λ , or by a matrix of the form

$$J_{\lambda, \bar{\lambda}} = \begin{pmatrix} R_{\lambda, \bar{\lambda}} & I & & & \\ & R_{\lambda, \bar{\lambda}} & I & & \\ & & \ddots & \ddots & \\ & & & R_{\lambda, \bar{\lambda}} & I \\ & & & & R_{\lambda, \bar{\lambda}} \end{pmatrix} \quad (4.11)$$

with

$$R_{\lambda, \bar{\lambda}} = \alpha I + \Omega = \begin{pmatrix} \alpha & \omega \\ -\omega & \alpha \end{pmatrix} \quad \text{and} \quad I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

for some couple of complex conjugated eigenvalues $\lambda = \alpha + i\omega$ and $\bar{\lambda} = \alpha - i\omega$, respectively.

Proof. According to the Jordan normal form theorem 4.5, there exists a basis of \mathbb{C}^N such that the complexified operator $L^{\mathbb{C}}$ is represented by a block diagonal Jordan matrix.

Consider a Jordan block with real eigenvalue λ . If $z = x + iy$ is a L_{λ} -cyclic vector, then either x or y is a real L_{λ} -cyclic vector. This real cyclic vector generates, therefore, a real Jordan chain of the same real dimension as the complex dimension of original block.

We now consider a Jordan block with complex eigenvalue $\lambda = \alpha + i\omega$. If λ is not real, then the complexified operator also admits an eigenvalue $\bar{\lambda} = \alpha - i\omega$, and a corresponding Jordan block of equal dimension. Indeed, if v is L_{λ} -cyclic, then \bar{v} is $L_{\bar{\lambda}}$ -cyclic. Proceeding as in the two dimensional case, one easily sees that this couple of complex Jordan blocks give origin to a real Jordan block of the form (4.11). \square

The invariant subspaces E_{λ} or $E_{\lambda, \bar{\lambda}}$ of theorem 4.6 are also referred to as *root spaces*, using a terminology borrowed from the theory of Lie algebras.

4.5 Hyperbolic linear flows

Stable and unstable spaces. Given a linear vector field L on \mathbb{R}^N , defined in the canonical basis by a real matrix A , we are interested in its flow $\Phi_t = e^{tL}$, which solves the linear system

$$\dot{x} = Ax.$$

We already saw that the asymptotic behavior of the flow e^{tL} in each root space depends on the sign of the real part of the corresponding eigenvalue.

One can write the total space as a direct sum of three invariant subspaces

$$\mathbb{R}^N = E^{-} \oplus E^0 \oplus E^{+}$$

where the *stable space* E^{-} is the direct sum of those root spaces with $\Re(\lambda) < 0$, the *unstable space* E^{+} is the direct sum of those root spaces with $\Re(\lambda) > 0$, and finally the *neutral space* E^0 is the direct sum of those root spaces with $\Re(\lambda) = 0$.

Sinks and sources. The linear system, or better its equilibrium point 0, is called a *sink* if all the eigenvalues have negative real part, i.e. $\Re(\lambda) < 0$, so that that $\mathbb{R}^N = E^{-}$. It is called a *source* if all the eigenvalues have positive real part, so that $\mathbb{R}^N = E^{+}$. It is clear that reversing the arrow of time transforms a sink to a source, and vice-versa, since $(e^{tL})^{-1} = e^{-tL}$.

Theorem 4.7. *The linear system $\dot{x} = L(x)$ is a sink iff it satisfies one of the following equivalent conditions:*

- i) all the eigenvalues of L have negative real part,
- ii) all solutions decay $e^{tL}v \rightarrow 0$ when $t \rightarrow \infty$,
- iii) there exist a positive $\alpha > 0$ and a constant C such that for all $v \in \mathbb{R}^N$

$$\|e^{tL} v\| \leq C e^{-\alpha t} \|v\| \quad \text{for times } t \geq 0. \quad (4.12)$$

Proof. It is obvious that iii) \Rightarrow ii). It is also clear that ii) \Rightarrow i), because if some eigenvalue has $\Re(\lambda) \geq 0$, then one easily find, in the corresponding Jordan chain, a solution which does not decay to zero. Finally, to see that i) \Rightarrow iii), we note that this holds in each Jordan block according to theorem 4.4. But if we have norms in each subspace of a direct sum decomposition (as for example the restrictions of the Euclidean norm), we can define a norm on the total by space taking their maximum (or their sum, or the square root of the sum of theirs squares). With respect to this norm, we then have the inequality (4.12) for some $\alpha > 0$ strictly smaller than all the $|\Re(\lambda)|$'s and some maximal constant C . Again, by the equivalence of all norms, the same inequality holds w.r.t. to the any norm in \mathbb{R}^N , for some possibly different constant C . \square

Thus, all trajectories of a sink decay exponentially fast to the origin. Conversely, all trajectories of a source are exponentially stretched, i.e. satisfy an inequality like

$$\|e^{tL} v\| \geq C e^{\beta t} \|v\|$$

for some $\beta > 0$ and all $t \geq 0$, and therefore diverge exponentially fast as $t \rightarrow \infty$, provided the initial condition is not the equilibrium, i.e. $v \neq 0$. If the linear field has non-real eigenvalues, trajectories may decay or diverge along logarithmic spirals.

Hyperbolic linear flows. A linear vector field L is called *hyperbolic* if the spectrum of its complexification is disjoint from the imaginary axis, i.e. if all the eigenvalues λ , real or complex, have non-zero real part $\Re(\lambda) \neq 0$. The total space of a hyperbolic vector field therefore splits as a direct sum

$$\mathbb{R}^N = E^- \oplus E^+$$

of a stable and an unstable invariant subspace.

Of course, sinks and sources are hyperbolic, but the most interesting case is when both the stable and the unstable subspaces are not-empty. Reasoning as in the proof of theorem 4.7, one shows that

Theorem 4.8. *Let L be a hyperbolic linear field. The phase space is a direct sum of two invariant subspaces $\mathbb{R}^N = E^- \oplus E^+$, the stable and the unstable subspaces, and there exist positive constants $\alpha, \beta > 0$ and a constant C such that*

$$\|e^{tL} v\| \leq C e^{-\alpha t} \|v\| \quad \text{if } v \in E^- \text{ and } t \geq 0$$

and

$$\|e^{-tL} v\| \leq C e^{-\beta t} \|v\| \quad \text{if } v \in E^+ \text{ and } t \geq 0$$

Thus, the flow of a hyperbolic linear vector field contracts vectors in the stable space and stretch vectors in the unstable space. Indeed, the stable and the unstable subspaces E^\pm may be characterized/defined as the sets of those vectors satisfying $e^{\pm tL} v \rightarrow 0$ for $t \rightarrow \infty$, respectively. If both spaces are not empty, generic trajectories, not starting in $E^- \cup E^+$, diverge for $t \rightarrow \pm\infty$.

It turns out that the hyperbolic vector fields are precisely the structurally stable linear vector fields. This is the starting point of a large area of the modern theory of dynamical systems, called *hyperbolic theory*. Classical references are [HS74, PM78].

5 Numbers and dynamics

Another important source of interesting dynamics is, quite surprisingly, elementary number theory.

5.1 Decimal expansion and multiplication by ten

Decimal expansion. When children we learn to represent numbers as decimals, like

$$3.14159265358979323846264338327950288419716939937510 \dots$$

Of course, there is nothing special with the number 10, it is but the number of fingers in our hands. Any other integer $d \geq 2$ would work. Representing a non-negative (for simplicity) real number $x \in \mathbb{R}_+$ in base 10 means writing x as the sum of a convergent series

$$\begin{aligned} x &= "X_m \dots X_2 X_1 X_0 . x_1 x_2 x_3 \dots" \\ &:= X_m \cdot 10^m + \dots + X_2 \cdot 10^2 + X_1 \cdot 10 + X_0 + \frac{x_1}{10} + \frac{x_2}{10^2} + \frac{x_3}{10^3} + \dots \\ &= \sum_{n=0}^m X_n \cdot 10^n + \sum_{n=1}^{\infty} x_n \cdot 10^{-n} \end{aligned}$$

where $X_n, x_n \in \{0, 1, 2, \dots, 9\}$ and $m \geq 0$ (the series above is absolutely convergent because it is bounded by 9 times the geometric series $\sum_{n=1}^{\infty} (1/10)^n$).

The finite sum

$$[x] := \sum_{n=0}^m X_n \cdot 10^n \in \mathbb{Z}$$

is the *integral part* of x , the largest of those integers n such that $n \leq x$. The possibly infinite sum

$$\{x\} := 0.x_1 x_2 x_3 \dots = \sum_{n=1}^{\infty} x_n \cdot 10^{-n} \in [0, 1)$$

is the *fractional part* of x , the difference $\{x\} = x - [x]$. Consequently, $[x] + \{x\} = x$.

Some representations terminate, i.e. have $x_n = 0$ starting from some $n \geq N$, and some others are *recurring* (or *eventually periodic*), i.e. of the form

$$[x] + 0.x_1 x_2 \dots x_k \overline{a_1 a_2 \dots a_n} := [x] + 0.x_1 x_2 \dots x_k a_1 a_2 \dots a_n a_1 a_2 \dots a_n a_1 a_2 \dots a_n \dots$$

for some finite recurring word $a_1 a_2 \dots a_n$ (and of course a terminating decimal is a recurring one with recurring word $\overline{0}$).

The representation is unique, hence defines a bijection between \mathbb{R} and the space of infinite words $X_m \dots X_2 x_1 X_0 . x_1 x_2 x_3 \dots$ as above, if we do not admit recurrent 9's, i.e. if we substitute $\dots x_{k-1} \overline{9}$ with $\dots (x_{k-1} + 1) \overline{0}$ (where we assume $x_{k-1} \neq 9$).

Division algorithm. The iterative scheme to obtain the decimal representation of a rational number is the “division algorithm” that we also learn when children. Consider a positive rational number $x = p/q$ with $p, q \in \mathbb{N}$:

$$\frac{p}{q} = x_0 . x_1 x_2 x_3 \dots$$

The integer $[x] = x_0$ is “the number of times q is contained in p ”, i.e. the unique integer such that

$$p = x_0 \cdot q + r_0$$

for some rest r_0 which is an integer $0 \leq r_0 < q$. Hence, $p/q = x_0 + r_0/q$ and $0 \leq r_0/q < 1$. The “geometric” meaning of x_1 is that the point r_0/q lies between $0.x_1$ and $0.x_1 + 0.1$. Multiplying by 10 and then by q this means that

$$x_1 \cdot q \leq 10 \cdot r_0 < x_1 \cdot q + q$$

or, equivalently, that x_1 is the unique integer between 0 and 9 such that

$$10 \cdot r_0 = x_1 \cdot q + r_1$$

where, again, the rest r_1 is a non-negative integer $0 \leq r_1 < q$. And so on. Hence, the digits of the decimal expansion of p/q are iteratively determined by

$$10 \cdot r_{n-1} = x_n \cdot q + r_n \quad \text{where} \quad 0 \leq r_n < q$$

Since the possibilities for the rests are finite, they necessarily recur. On the other side, a simple computation shows that a recurring decimal is a (series converging to a) rational number.

Theorem 5.1. *The rational numbers are precisely those real numbers whose representation in base 10 (or any other base $d \geq 2$) is (eventually) repeating/recurring.*

Meanwhile, there exist irrational numbers. For example,

$$0.101001000100001 \dots = \frac{1}{10} + \frac{1}{10^3} + \frac{1}{10^6} + \frac{1}{10^{10}} + \frac{1}{10^{15}} + \dots$$

is irrational, since it is not recurring.

Indeed, almost all numbers are irrational, in a precise probabilistic sense, since rationals are countable.

The weight of the rationals. Consider the unit interval $I = [0, 1]$, and imagine to cut out all its rational points. What is left is a set, $I \setminus \mathbb{Q}$, whose length is equal to the length of the interval! Indeed, the rationals are countable, for example those inside I may be ordered according to

$$0 \quad 1 \quad 1/2 \quad 1/3 \quad 2/3 \quad 1/4 \quad 3/4 \quad 1/5 \quad 2/5 \quad 3/5 \quad 4/5 \quad \dots$$

say $I \cap \mathbb{Q} = \{r_1, r_2, r_3, \dots\}$. Given an (arbitrarily small) $\varepsilon > 0$, we may even cut out a whole interval $J_n = (r_n - \ell_n/2, r_n + \ell_n/2)$ of finite diameter $\ell_n = \varepsilon/2^n$ around each r_n . The measure of what is left of the unit interval is

$$\text{length} (I \setminus (\cup_{n=1}^{\infty} J_n)) \geq 1 - \sum_{n=1}^{\infty} \varepsilon/2^n = 1 - \varepsilon$$

In other words, the rationals inside the unit interval have neighbourhoods of arbitrarily small length! Mathematicians say that

Theorem 5.2. *Rationals form a set of Lebesgue measure zero inside the real line.*

Therefore, almost all numbers are irrationals. In other words, if we “choose” a random number in the interval $[0, 1]$, with respect to the uniform distribution giving probability $|b - a|$ to any interval $[a, b] \subset [0, 1]$, it will be irrational “with probability one”. Despite this fact, showing that a “given” real number, like π , e , ... is irrational may be a hard problem.

ex: Show that the decimal representation of a (reduced) rational p/q terminates iff the denominator is of the form $q = 2^\alpha 5^\beta$ for some non-negative integers α and β .

ex: Write $1/3$ in base 2, and $2/3$ in base 3 and 7.

ex: Show that the decimal (or any other base) representation of a rational number is repeating (observe that the possibilities for the rests r_n are finite). Then show the converse: a repeating decimal represents a rational number (compute the sum of the series).

ex: Give examples of non-repeating decimal expansions (see [HW59], section 9.4).

ex: Prove that *Euler's number*

$$e := \sum_{n=0}^{\infty} \frac{1}{n!}$$

is irrational (Fourier's idea: assume that $e = p/q$ for some positive integers p and q , and deduce that $x = q!(e - \sum_{n=0}^q 1/n!)$ is then an integer. Estimate the series $x = \sum_{n=q+1}^{\infty} q!/n!$ and prove that $0 < x < 1$).

Multiplication by an integer. The representation of real number in base 10 is strictly related to the dynamics of a particular transformation acting on the circle.

Let $N \geq 2$ be an integer. The transformation $F_N : \mathbb{R} \rightarrow \mathbb{R}$ sending each number x to its multiple $F_N(x) = Nx$ has a trivial dynamics, since all trajectories diverge, apart from the fixed point 0. Things get interesting if we do not allow trajectories to escape, i.e. if we force them to a bounded domain. One way to do it is considering the quotient circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$, and define the transformation $E_N : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ as

$$x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$$

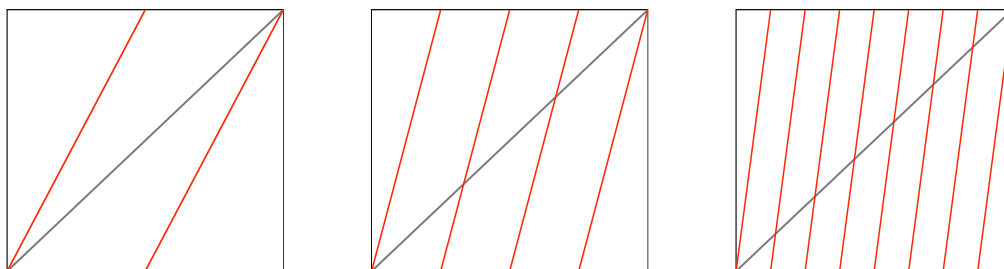
Let $\pi : \mathbb{R} \rightarrow \mathbb{R}/\mathbb{Z}$ denotes the projection of the real line over the circle, sending $\pi(x) = x + \mathbb{Z}$. Then $\pi \circ F_N = E_N \circ \pi$, i.e., F_N is a lift of E_N .

Alternatively, we could have defined a transformation of the unit interval $[0, 1]$ into itself sending $x \mapsto \{Nx\} = Nx - [Nx]$, thus avoiding the identification $0 \sim 1$.

A “fundamental domain” for the action of \mathbb{Z} on the real line is the interval $[0, 1)$. This means that any class $x + \mathbb{Z} \in \mathbb{R}/\mathbb{Z}$ admits one and only one representative $x \in [0, 1)$. If $x = 0.x_1x_2x_3 \dots$ is the representation of $x \in \mathbb{R}/\mathbb{Z} \simeq [0, 1)$ in base N , then E_N sends

$$0.x_1x_2x_3 \dots \mapsto 0.x_2x_3x_4 \dots$$

The simplest case is that of the *doubling map*, $E_2(x + \mathbb{Z}) = 2x + \mathbb{Z}$.



Graph of the doubling map, and its first two iterates.

ex: Find the cardinality of the inverse image by E_N of a generic point in \mathbb{R}/\mathbb{Z} .

ex: Find periodic and pre-periodic points of E_N , show that they are dense in the circle.

ex: Show that the identification $h : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{S}^1$, given by $x + \mathbb{Z} \mapsto e^{2\pi i x}$, is a topological conjugation between the doubling map E_2 and the restriction of the squaring map $z \mapsto z^2$ to the unit circle $\mathbb{S} \subset \mathbb{C}$. State the corresponding result for the multiplication by an arbitrary integer $N \geq 2$.

Multiplication by 10 and frequencies of digits. Consider a (random?) number $x \in [0, 1)$, and its decimal representation $x = 0.x_1x_2x_3 \dots$, with digits $x_k \in \{0, 1, 2, \dots, 9\}$ and no recurring 9's. The first digit x_1 is, for example, equal to 7 iff x belongs to the interval $A = [0.7, 0.8) + \mathbb{Z} \subset \mathbb{R}/\mathbb{Z}$ (though as an interval of the circle). Since $E_{10}(0.x_1x_2x_3 \dots + \mathbb{Z}) = 0.x_2x_3x_4 \dots + \mathbb{Z}$, the second digit x_2 is equal to 7 iff $E_{10}(x + \mathbb{Z}) \in A$. And so on. Let $\chi_A : \mathbb{R}/\mathbb{Z} \rightarrow \{0, 1\}$ be the characteristic

function of the interval A , which is equal to $\chi_A(x + \mathbb{Z}) = 1$ iff $x + \mathbb{Z}$ belongs to A . There follows from the above discussion that the frequency of 7's between the first n digits of x is a Birkhoff sum

$$\frac{1}{n} \text{card}\{1 \leq k \leq n : x_k = 7\} = \frac{1}{n} \sum_{k=0}^{n-1} \chi_A(E_{10}^k(x + \mathbb{Z}))$$

The limit for $n \rightarrow \infty$, if it exists, should be interpreted as an “asymptotic frequency” of 7's in the decimal representation of x .

It is clear that similar considerations holds for other possible values of x_1 , say 3 or 9, as well as for possible initial finite strings $x_1x_2 \dots x_k$, as 123 or 1415. Therefore, the dynamics of E_{10} on the circle contains informations on the patterns of decimal representations of numbers.

ex: Conjecture a value for the asymptotic frequency of 7's, for “typical” numbers x .

ex: Show that any value, for example a rational $0 \leq p/q \leq 1$, may be attained as an asymptotic frequency of 7's, for particular numbers x .

ex: Show that there exist numbers for which the limit does not exist.

5.2 Bernoulli shifts

The abstract version of multiplication by N on the circle is the shift on the space of Bernoulli trials, a map which is basic in probability.

Infinite words. Let $\mathcal{A} \approx \{1, 2, \dots, N\}$ be an “alphabet” made of $N \geq 2$ letters, i.e. a finite set equipped with the discrete topology, and let $\Sigma^+ := \mathcal{A}^{\mathbb{N}}$ be the topological product of infinite copies of \mathcal{A} . Points of Σ^+ are actually sequences $(x_n)_{n \in \mathbb{N}}$ with values in \mathcal{A} , but are more conveniently denoted as

$$x = x_1x_2x_3 \dots x_n \dots$$

with $x_n \in \mathcal{A}$, and interpreted as “infinite words” in the letters of the alphabet \mathcal{A} . In probability theory, the x_k 's represents the outcomes of a sequence of trials of some experience with N possible outcomes (as a dice with N faces).

The *product topology* is the weakest topology on $\mathcal{A}^{\mathbb{N}}$ such that all the projections $\pi_n : \Sigma^+ \rightarrow \mathcal{A}$, sending $x \mapsto x_n$, are continuous. A basis for this topology is the family \mathcal{C} of centered cylinders. A *centered cylinder* is a subset

$$C_\alpha := \{x \in \Sigma^+ \text{ s.t. } x_1 = \alpha_1, x_2 = \alpha_2, \dots, x_k = \alpha_k\}$$

where $\alpha = \alpha_1\alpha_2 \dots \alpha_k \in \mathcal{A}^k$ is a finite word of length $k \in \mathbb{N}$. More colloquially, C_α is the set of those infinite words x starting with the finite word α , i.e. of the form $x = \alpha*$, with an obvious meaning of the symbol “*” (as in the UNIX language). Thus, a basis \mathcal{C} of the product topology is the countable family of C_α , when α ranges in the set $\bigcup_k \mathcal{A}^k$ of all finite words in the letters of \mathcal{A} . By definition, an open set of the topological product Σ^+ is a union $A = \bigcup_\alpha C_\alpha$ of centered cylinders.

Observe that the family of centered cylinders is a basis of a topology because it is covering, since obviously $\Sigma^+ = C_1 \cup C_2 \cup \dots \cup C_N$, and because the intersection of two cylinders is the empty set or one of the two cylinders. Indeed, two cylinders C_α and C_β have non-empty intersection iff one of the two words, say $\alpha = \alpha_1\alpha_2 \dots \alpha_k$, is the initial string of the other word, in the sense that $\beta = \alpha_1\alpha_2 \dots \alpha_k\beta_{k+1} \dots \beta_{k+i}$, and in this case $C_\alpha \cap C_\beta = C_\beta$. The idea is that the longer is the word α the smaller is the cylinder C_α .

Ultrametrics. The product topology is metrizable. This means that there exist metrics on Σ^+ which induce the product topology. One possibility is the metric

$$d_\lambda(x, y) = \sum_{n=1}^{\infty} \lambda^{-n} \cdot |x_n - y_n|$$

for some $\lambda > 1$ (for example $\lambda = N$). Another possibility, simpler to deal with, is to define $\text{ord}(x, y) := \min\{n \in \mathbb{N} : x_n \neq y_n\}$, the smallest place where the two words x and y differ, and then a distance as

$$d_\infty(x, y) = \lambda^{-\text{ord}(x, y)}$$

if $x \neq y$, and zero otherwise. It is clear that centered cylinders C_α are both closed and open balls for this metric, as strange as it may seem. It turns out that this is indeed an *ultrametric*, triangular inequality being a consequence of the stronger ultrametric inequality

$$d_\infty(x, y) \leq \max\{d_\infty(x, z), d_\infty(z, y)\}$$

Between the strange properties of ultrametric spaces, one verifies that any point of a ball is its center. The space Σ^+ is the abstract example of a Cantor set: a compact, perfect and totally disconnected metric space.

ex: Show that d_∞ is an ultrametric.

ex: Show that any point of a ball in a ultrametric space is its center, and that balls are both open and closed (you may find paradoxical a ball all of whose points are its centers; meanwhile, we know, at least since Eratosthenes (around 200 BC), that we are living in such a ball: what is it?).

Bernoulli shift. The *Bernoulli shift* is the transformation $\sigma : \Sigma^+ \rightarrow \Sigma^+$ which “forgets the first letter” of the infinite word, sending

$$x_1x_2x_3 \cdots \mapsto \sigma(x_1x_2x_3) := x_2x_3x_4 \cdots$$

It is continuous, because the inverse image of any centered cylinder is a union of centered cylinders, hence an open set. It is not invertible, and indeed the inverse image of any point is made of N different points (the choices for the first letter of the infinite word).

In probability, letters of the alphabet \mathcal{A} represent the possible outcomes of an experience, as tossing a coin or a dice. An infinite word $x_1x_2x_3 \cdots x_n \cdots$ therefore represents the successive results of a countable set of experiences, ordered, for example, by time n . Iteration of σ means forgetting the outcomes of the first experiences.

ex: Describe periodic and pre-periodic points of σ . Show that they are dense in Σ^+ .

ex: Consider the alphabet $\mathcal{A} = \{0, 1, 2, \dots, 9\}$. Define a map $h : \mathcal{A}^{\mathbb{N}} \rightarrow [0, 1]$ as

$$x_1x_2x_3 \cdots \mapsto 0.x_1x_2x_3 \cdots$$

Show that h is a semi-conjugation between the shift σ and the multiplication E_{10} on the circle.

5.3 Rotations of the torus

Rotations of the circle. The *circle*, or *one-dimensional torus*, is the quotient $\mathbb{T} := \mathbb{R}/\mathbb{Z}$ of the commutative group \mathbb{R} modulo its subgroup \mathbb{Z} , equipped with the quotient topology. We denote by $\pi : \mathbb{R} \rightarrow \mathbb{T}$ the projection $x \mapsto \pi(x) := x + \mathbb{Z}$. The euclidean metric on the real line induces a metric on the circle, defined by

$$\begin{aligned} d(x + \mathbb{Z}, y + \mathbb{Z}) &= \min_{x' \in x + \mathbb{Z}, y' \in y + \mathbb{Z}} |x' - y'| \\ &= \min_{n \in \mathbb{Z}} |x - y + n| \end{aligned}$$

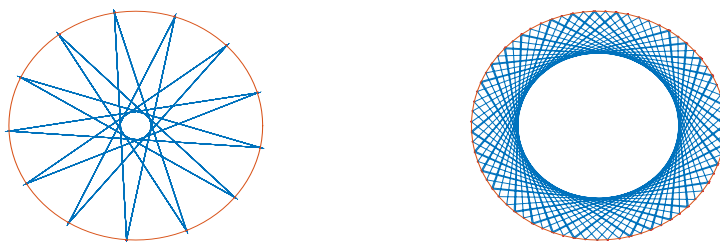
Thus, the distances between the classes of x and y is the minimal distance between the subsets $x + \mathbb{Z}$ and $y + \mathbb{Z}$ of the real line. Observe that the diameter of the circle, i.e. the maximal distance between two points, is $1/2$.

Rotations of the circle are the transformations $R_\alpha : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ defined by

$$x + \mathbb{Z} \mapsto x + \alpha + \mathbb{Z}$$

where $\alpha \in \mathbb{R}$. Observe that \mathbb{R}/\mathbb{Z} is a commutative group, and the R_α are its translations, since only the class $\alpha + \mathbb{Z}$ of α matters. Also, rotations are the isometries of the circle which preserve the orientation. If we identify the circle with the unit circle $\mathbb{S}^1 := \{z \in \mathbb{C} \text{ t.q. } |z| = 1\} \subset \mathbb{C}$ in the complex plane, by means of the homeomorphism $x + \mathbb{Z} \mapsto e^{2\pi i x}$, rotations are the transformations $z \mapsto e^{i2\pi\alpha} z$.

It is interesting to observe that trajectories of a circle rotations are the successive points where a billiard ball hits the boundary circle if thrown inside a circular billiard.



Orbits of a circular billard, with rational and irrational angle.

Theorem 5.3. *A rotation R_α has periodic points iff α is rational.*

Proof. If α is rational, and equal to the reduced fraction p/q , then all points are periodic with period q , since $x + q\alpha + \mathbb{Z} = x + \mathbb{Z}$ for any x . On the other side, if α is irrational, there exists no natural $n \geq 1$ such that $x + \mathbb{Z} = x + n\alpha + \mathbb{Z}$, independently on x . \square

Indeed, what we showed is that all orbits of a rational rotation are periodic, hence finite. On the other side, all orbits of an irrational rotation are infinite, and this will be the interesting case.

ex: Observe that any point $x + \mathbb{Z}$ of the circle \mathbb{R}/\mathbb{Z} has a unique representative $\{x\}$ in the interval $[0, 1)$, the fractional part of x , and show that the distance between two points of the circle is given by the explicit formula

$$d(x + \mathbb{Z}, y + \mathbb{Z}) = \min\{|\{x\} - \{y\}|, 1 - |\{x\} - \{y\}|\}$$

Rotations of a torus. The N -dimensional *torus* is the quotient $\mathbb{T}^N := \mathbb{R}^N/\mathbb{Z}^N$, equipped with the quotient metric. *Rotations* are the homeomorphisms $R_\alpha : \mathbb{T}^N \rightarrow \mathbb{T}^N$ defined by

$$x + \mathbb{Z}^N \mapsto x + \alpha + \mathbb{Z}^N$$

where now $\alpha \in \mathbb{R}^N$.

ex: Try to understand possible orbits of a torus rotation.

5.4 Dyadic adding machine

p -adic number fields and integers. The field \mathbb{R} of real numbers may be considered (actually constructed) as the completion of the rational number field \mathbb{Q} with respect to the Euclidean norm $|x|_\infty := \max\{\pm x\}$. This means that real numbers are equivalence classes of fundamental sequences of rationals, two fundamental sequences (x_n) and (y_n) being in the same class, i.e. representing the same real ‘ $x := \lim_{n \rightarrow \infty} x_n$ ’, if $|x_n - y_n|_\infty \rightarrow 0$.

It happens that there exist other “norms”, i.e. positive and homogeneous functionals $x \mapsto |x|$ on the rationals satisfying the triangular inequality, which respect the multiplicative structure, i.e. such that $|xy| = |x||y|$. Such norms are called *valuations*.

Let $p = 2, 3, 5, 7, \dots$ be a rational prime, also called *place* in this context. The *order* of a non-zero rational $x \in \mathbb{Q}^\times := \mathbb{Q} \setminus \{0\}$ at the place p is the unique integer $\text{ord}_p(x) = n$ such that $x = p^n a/b$ for some $a, b \in \mathbb{Z}$ which are not divided by p . The *p -adic valuation/place* is the absolute value on \mathbb{Q} defined as

$$|x|_p := p^{-\text{ord}_p(x)} \quad \text{if } x \neq 0$$

and $|0|_p = 0$. Clearly $\text{ord}_p(xy) = \text{ord}_p(x) + \text{ord}_p(y)$ (just like the degree of polynomials), and this gives homogeneity of $|\cdot|_p$. Triangular inequality follows from the observation that

$$\text{ord}_p(x + y) \geq \min\{\text{ord}_p(x), \text{ord}_p(y)\}$$

One can show that those, together with the euclidean norm, are the only valuations on \mathbb{Q} , modulo trivial equivalences. The *p -adic (topological) number field* \mathbb{Q}_p is the completion of \mathbb{Q} with respect to $|\cdot|_p$ (uniqueness is trivial, and existence may be proved as usual considering equivalence classes of fundamental sequences, the only annoying issue being keeping track of the field operations). The *p -adic valuation*, naturally extended to \mathbb{Q}_p , is “non-Archimedean” (i.e. does not satisfy the “Archimedean property” that for all $\varepsilon > 0$ and all N there exists an integer n such that $n\varepsilon > N$) since triangular inequality is enhanced by the stronger “ultra-metric” inequality

$$|x + y|_p \leq \max\{|x|_p, |y|_p\}.$$

This causes many paradoxical properties. For example, closed balls are open as well (hence called “clopen”), and any point of a ball is its center. The ultrametric inequality also implies that

$$|b_1 + b_2 + \dots + b_n|_p \leq \max_{1 \leq k \leq n} |b_k|_p.$$

Consequently, a series $\sum_{n=0}^{\infty} b_n$ converges (for the p -adic metric, of course!) iff the norm of its terms $|b_n|_p \rightarrow 0$ as $n \rightarrow \infty$ (there is no room for divergent series like the harmonic series in the p -adic world!).

The ring of *p -adic integers* \mathbb{Z}_p is the closure of \mathbb{Z} in \mathbb{Q}_p . One can describe the p -adic integers as the inductive limit $\mathbb{Z}_p = \varprojlim (\mathbb{Z}/p^n\mathbb{Z})$, and represent a p -adic integer as a series

$$z = \dots z_n \dots z_2 z_1 z_0 := \sum_{n=0}^{\infty} z_n p^n \quad (5.1)$$

with $z_n \in \mathcal{A}_p := \{0, 1, 2, \dots, p-1\} \approx \mathbb{Z}/p\mathbb{Z}$, which converges in \mathbb{Q}_p because the norm of the generic term $z_n p^n$ is bounded by $|z_n p^n|_p = p^{-n} \rightarrow 0$ as $n \rightarrow \infty$. Thus, as a topological space (not as a ring!), \mathbb{Z}_p is isomorphic to the topological product $\Sigma^+ = \mathcal{A}^{\mathbb{N}}$, the space of infinite words (written backwards!) $\dots z_n \dots z_3 z_2 z_1$ in the letters of the alphabet \mathcal{A}_p . Observe that $\mathbb{Z}_p = \{x \in \mathbb{Q}_p \text{ s.t. } |x|_p \leq 1\}$, i.e. the ring of p -adic integers is the clopen ball of radius one around 0 in \mathbb{Q}_p .

Any p -adic number $x \in \mathbb{Q}_p$ can be represented uniquely as $x = z + r$ where $z = [x]_p = \sum_{n \geq 0} x_n p^n \in \mathbb{Z}_p$ is the “ p -adic integer part” and $r = \{x\}_p = \sum_{n=1}^{\infty} x_n p^{-n} \in \mathbb{Z}[1/p]$ is the “ p -adic fractional part”. In symbols,

$$x = \dots x_n \dots x_2 x_1 x_0 x_{-1} \dots x_{-N} = \sum_{n=-N}^{\infty} x_n p^n$$

The quotient $\mathbb{Q}_p/\mathbb{Z}_p = \mathbb{Z}[1/p]/\mathbb{Z}$ is a discrete (additive) group where the norm $|\cdot|_p$ takes values p^n with $n \in \mathbb{N}$. Multiplication by p is a uniform contraction $x \mapsto px$ of \mathbb{Q}_p with Lipschitz constant p^{-1} , and its inverse $x \mapsto xp^{-1}$ uniformly expands distances by a factor p . Thus, \mathbb{Z}_p is the disjoint union $\cup_{a_i=0}^{p-1} (a_i + p\mathbb{Z}_p)$ of p clopen balls of radius p^{-1} (and so on, iterating the contraction). Also, one can represent the field of p -adic numbers as a union $\mathbb{Q}_p = \bigcup_{n \in \mathbb{N}} p^{-n}\mathbb{Z}_p$.

ex: Compute the following sums in \mathbb{Z}_2 .

$$\dots 011 + \dots 001 \quad \dots 0101 + \dots 1010 \quad \dots 111 + \dots 001$$

Find the additive opposite of $1 = \dots 001$ in \mathbb{Z}_2 .

Adding machine. Consider the ring \mathbb{Z}_2 of dyadic integers, thought as an additive topological group. The *dyadic adding machine* (or “*Kakutani-von Neumann odometer*”) is the translation $f : \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$, defined by

$$z \mapsto z + 1$$

Observe that f changes the first (starting from the right) digit of $x = \dots x_2 x_1 x_0$, and consequently its n -th iterate changes the n -th digit of x . Thus, we have one more example of a translation in a compact topological group without periodic points.

ex: Show that $z \mapsto z + 1$ is a homeomorphism of \mathbb{Z}_2 , and find its inverse.

5.5 Continued fractions and Gauss map

Continued fractions. Any real number $x \in \mathbb{R}$ can be represented (uniquely if irrational, and with only a minor ambiguity if rational) as a *continued fraction*

$$x \sim [a_0; a_1, a_2, a_3, \dots] := a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \frac{1}{\ddots}}}}$$

with $a_0 \in \mathbb{Z}$ and “partial quotients” $a_n \in \mathbb{N}$ if $n \geq 1$. This means that x is equal to the limit of the *convergents*, the finite continued fractions defined as

$$p_n/q_n = [a_0; a_1, a_2, \dots, a_n] := a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\ddots + \frac{1}{a_n}}}}$$

as $n \rightarrow \infty$. Observe that finite continued fractions are rationals, and obey the recursion

$$[a_0; a_1, a_2, \dots, a_n, a_{n+1}] = [a_0; a_1, a_2, \dots, a_n + 1/a_{n+1}] \quad (5.2)$$

Continued fractions constitute the fundamental tool to investigate rational approximations to real numbers, because they provide base-free, hence intrinsic, rational approximations. Thus, while Earthlings with ten fingers write $\pi = 3.1415\dots$ and Martians with three fingers write $\pi = 10.0102\dots$, they all agree to write $\pi = [3; 7, 15, 1, 292, \dots]$. Moreover, they provide the best rational approximations, in a certain precise sense [HW59, Kh35].

Construction and Gauss map. The continued fraction converging to a given number $x \in \mathbb{R}$ is given essentially by Euclid's algorithm to find the m.c.d. of two integers. One starts with $a_0 = \lfloor x \rfloor \in \mathbb{Z}$ (here the "floor" function $\lfloor x \rfloor$ returns the smallest integer n such that $n \leq x < n+1$), and write $x = a_0 + x_0$ for some $x_0 = \{x\} \in [0, 1)$. Then define the *Gauss map* $G : (0, 1] \rightarrow [0, 1]$ as

$$G(x) := 1/x - \lfloor 1/x \rfloor, \quad (5.3)$$

(thus, $G(x)$ is the fractional part of the inverse of x) and inductively define the partial quotients $a_n \in \mathbb{N}$ and the "rests" $x_n \in [0, 1)$ as

$$a_{n+1} = \lfloor 1/x_n \rfloor \quad x_{n+1} = G(x_n),$$

provided all the $x_0, x_1, \dots, x_n \neq 0$. Then,

$$x = a_0 + x_0 = a_0 + \frac{1}{a_1 + x_1} = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + x_2}} = \dots = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\ddots + \frac{1}{a_n + x_n}}}}$$

If some $x_n = 0$, the iteration stops and x is equal to a finite continued fraction as above. Conversely, if $x = p/q$ is rational, all the x_n 's are positive rationals, and have strictly decreasing denominators (for if $x_n = a/b$, then $1/x_{n+1} = x_n - a_n = (a - a_n b)/b = c/b$, and $c < b$ because $x_n - a_n < 1$). So, there must be some first x_n which is an integer, and the algorithm stops. Thus, finite continued fractions correspond to rationals (and are unique if we demand the last non-zero partial quotient be $a_n > 1$).

Convergence of the convergents. We must therefore understand the case of infinite continued fractions, which, as we already know, correspond to irrationals. The key observation is that the convergents $p_n/q_n = [a_0; a_1, a_2, \dots, a_n]$ of a continued fraction $[a_0, a_1, a_2, a_3, \dots]$ are determined by the partial quotients a_n 's according to the following recursive equation.

Theorem 5.4. *The convergents $p_n/q_n = [a_0; a_1, a_2, \dots, a_n]$ are obtained from the coefficients a_k 's by the recursions*

$$\begin{aligned} p_n &= a_n p_{n-1} + p_{n-2} \\ q_n &= a_n q_{n-1} + q_{n-2} \end{aligned} \quad (5.4)$$

given the initial conditions $p_0 = a_0$, $q_0 = 1$, and $p_{-1} = 1$, $q_{-1} = 0$ (or $p_{-2} = 0$ and $q_{-2} = 1$).

Proof. The proof is by induction. The first two values are easily verified. Assume the results holds until n , and compute

$$\begin{aligned} \frac{p_{n+1}}{q_{n+1}} &= [a_0; a_1, a_2, \dots, a_n, a_{n+1}] \\ &= [a_0; a_1, a_2, \dots, a_n + 1/a_{n+1}] \\ &= \frac{(a_n + 1/a_{n+1})p_{n-1} + p_{n-2}}{(a_n + 1/a_{n+1})q_{n-1} + q_{n-2}} \\ &= \frac{a_{n+1}(a_n p_{n-1} + p_{n-2}) + p_{n-1}}{a_{n+1}(a_n q_{n-1} + q_{n-2}) + q_{n-1}} \\ &= \frac{a_{n+1}p_n + p_{n-1}}{a_{n+1}q_n + q_{n-1}} \end{aligned}$$

□

It is important to write the recursion (5.4) in matrix notation as

$$\begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} = \begin{pmatrix} p_{n-1} & p_{n-2} \\ q_{n-1} & q_{n-2} \end{pmatrix} \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix},$$

whose solution, taking care of the initial conditions, is the backward product

$$\begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} = \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_1 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix} \quad (5.5)$$

of $n + 1$ integer matrices with determinant -1 . In particular,

$$p_n q_{n-1} - p_{n-1} q_n = (-1)^{n+1},$$

which says that the matrix with columns p_n, q_n and p_{n-1}, q_{n-1} is unimodular, i.e. belongs to the group $\text{GL}_2(\mathbb{Z})$ of (invertible) integer matrices with determinant ± 1 (two by two matrices whose rows and columns are relatively prime integers, a group which contains much arithmetical information!). This shows that the fractions p_n/q_n obtained using the recurrence in theorem 5.4 are reduced.

There also follows that

$$\frac{p_{n+1}}{q_{n+1}} - \frac{p_n}{q_n} = \frac{(-1)^n}{q_{n+1}q_n}$$

Also, one can easily show that convergents with n even form an increasing sequence, and convergents with n odd form a decreasing sequence. Another consequence of the recursion (5.4) is that the denominators q_n of an infinite continued fraction (i.e. such that $a_n \geq 1$ for all $n \geq 1$) satisfy

$$q_{n+2} \geq q_{n+1} + q_n \geq 2q_n,$$

and therefore grow exponentially fast:

$$q_{n+1} \geq 2^{n/2}.$$

(indeed, they grow at least like the Fibonacci sequence starting with $f_0 = 1$ and $f_1 = 1$, hence like $q_n \geq c\phi^n$, where $\phi = (1 + \sqrt{2})/2$ is the “ratio” and c is some positive constant). This implies that

$$\left| \frac{p_{n+1}}{q_{n+1}} - \frac{p_n}{q_n} \right| = \frac{1}{q_{n+1}q_n} \leq \frac{2}{2^n}$$

and therefore the sequence of the convergents p_n/q_n is fundamental. Its limit $\lim_{n \rightarrow \infty} p_n/q_n = x$ is called “value” of the continued fraction, and denoted by

$$x = [a_0; a_1, a_2, a_3, \dots] := \lim_{n \rightarrow \infty} [a_0; a_1, a_2, \dots, a_n].$$

It is also possible to find a lower bound to the difference between an irrational x and its convergents, and the two-sided estimate reads as follows:

$$\frac{1}{q_n(q_{n+1} + q_n)} < \left| x - \frac{p_n}{q_n} \right| < \frac{1}{q_{n+1}q_n}$$

ex: Use the quadratic equation $\phi^2 - \phi - 1 = 0$ to show that the “ratio” ϕ has the simplest continued fraction, namely

$$\frac{1 + \sqrt{5}}{2} = [1; 1, 1, 1, 1, \dots]$$

(observe that $\phi^{-1} = \phi - 1$ is a root of $x^2 + x - 1 = 0$, hence $x = 1/(1 + 1/x)$, and so on). Its convergents are 1 , 2 , $3/2$, $5/3$, $8/5$, $13/8$, $21/13$, $34/21$, \dots , ratios between successive Fibonacci numbers. It is also the (irrational) number with worse rational approximations, namely $|\phi - p/q| > (1/\sqrt{5})/q^2$ for any rational p/q .

ex: Also, the most famous irrational has a simple continued fraction. Show that

$$\sqrt{2} = [1; 2, 2, 2, 2, 2, \dots]$$

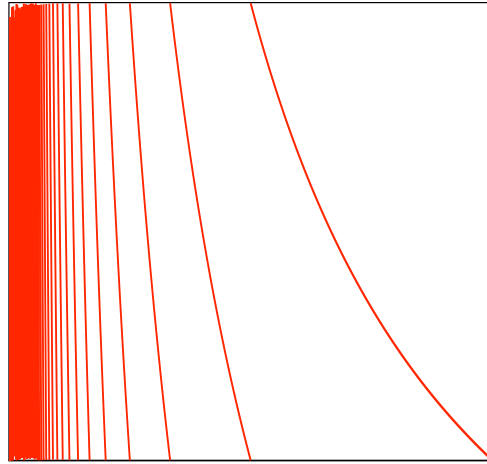
(observe that $1 + \sqrt{2}$ is the positive root of $x^2 - 2x - 1$. Hence $x = 2 + 1/x$, and so on). Its convergents are 1 , $3/2$, $7/5$, $17/12$, $41/29$, $99/70$, $239/169$, $577/408$, \dots

Continued fractions and Bernoulli shift. The continued fraction development, the map

$$x \mapsto [a_0; a_1, a_2, \dots]$$

realizes a conjugation between the restriction of the Gauss map (5.3) to the irrationals, the transformation $G : (0, 1] \setminus \mathbb{Q} \rightarrow (0, 1] \setminus \mathbb{Q}$, and the shift $\sigma : \mathbb{N}^{\mathbb{N}} \rightarrow \mathbb{N}^{\mathbb{N}}$ over an alphabeth $\mathcal{A} = \mathbb{N}$ of infinite letters. Indeed,

$$G : [0; a_1, a_2, a_3, \dots] \mapsto [0; a_2, a_3, a_4, \dots] .$$



Graph of the Gauss map.

ex: Find the (largest two or three) fixed points of the Gauss map, and compute their values.

Periodic continued fractions and quadratic irrationals. *Quadratic irrationals* (or *quadratic surds*) are irrational roots of quadratic polynomials with integer coefficients

$$f(x) = ax^2 + bx + c$$

(with $a, b, c \in \mathbb{Z}$), i.e. numbers like

$$x = \frac{\alpha + \sqrt{\beta}}{\delta}$$

where $\alpha, \beta, \delta \in \mathbb{Z}$, $\delta \neq 0$ and $\beta > 0$ which is not a square.

Theorem 5.5 (Lagrange). *The continued fraction of an irrational number $x \in \mathbb{R} \setminus \mathbb{Q}$ is periodic iff x is a quadratic irrational.*

See [Kh35, HW59].

5.6 Exponential sums

Arithmetic progressions . The dynamics of an *arithmetic progression*

$$a \quad a + \alpha \quad a + 2\alpha \quad a + 3\alpha \quad \dots \quad a + n\alpha \quad \dots,$$

obtained from the initial condition $x_0 = a$ using the recursion $x_{n+1} = x_n + \alpha$, is quite trivial. All trajectories $x_n = a + n\alpha$ diverge, provided $\alpha \neq 0$.

Something interesting happens if we compute time averages of the basic character of the real line, the observable $e : \mathbb{R} \rightarrow \mathbb{S} \subset \mathbb{C}$ given by

$$e(x) := e^{2\pi i x}.$$

Apart from a constant factor $e^{2\pi i a}$ and the normalization $1/N$, the Birkhoff averages of an arithmetic progression are

$$S_N(\alpha) = \sum_{n=0}^{N-1} e^{2\pi i \alpha n}.$$

ex: Show that the sum of the first n terms of an arithmetic progression $x_k = a + k\alpha$ is

$$\sum_{k=0}^{n-1} x_k = \frac{n}{2}(x_0 + x_{n-1}) = na + \frac{n(n-1)}{2}\alpha$$

Exponential sums. Sums as

$$E(N) = \sum_{n=1}^N e^{2\pi i x_n}$$

are called *exponential sums*, and contain “spectral information” about the distribution of the sequence of numbers (x_n) modulo 1. Triangular inequality gives the trivial bound $|E(N)| \leq N$, i.e. $E(N) = \mathcal{O}(N)$. If the different exponentials $e^{2\pi i x_n}$ were “uncorrelated”, as successive positions of a random walk in the plane, we should expect $E(N) = \mathcal{O}(\sqrt{N})$. This, of course, does not happen with “deterministic” generic sequences. The best we can hope is some bound as $E(N) = o(N)$ (which, in our case, would mean that the Birkhoff averages $\bar{\varphi}_n \rightarrow 0$).

ex: Observe that, for integer $q \geq 1$, the complex number $z = e^{2\pi i/q}$ is a non-trivial q -th root of unity. Hence,

$$1 + z + z^2 + \dots + z^{q-1} = 0.$$

Deduce that if $\alpha = p/q \in \mathbb{Q}$ with $p \in \mathbb{Z}$, then

$$\sum_{n=0}^{q-1} e^{2\pi i (p/q)n} = 1$$

so that the exponential sum $S_n(p/q)$ is periodic, and in particular is $\mathcal{O}(1)$.

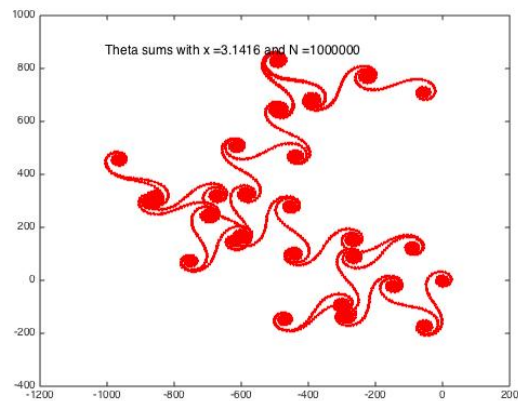
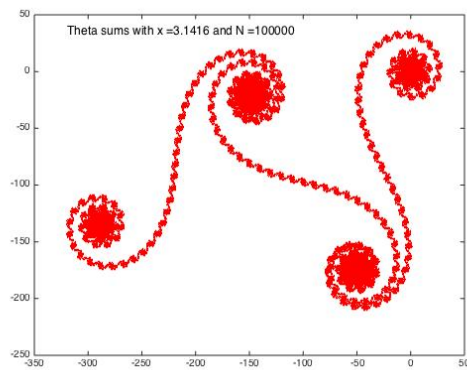
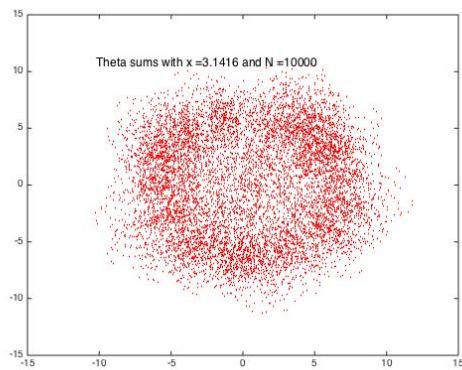
Gauss sums. Much more interesting are exponential sums defined by a “quadratic progression” $x_n = \alpha n^2$. These are

$$G_N(\alpha) = \sum_{n=0}^{N-1} e^{2\pi i \alpha n^2}.$$

When $\alpha = p/q$ is a rational, they are called (*quadratic*) *Gauss sums*, and they are extremely interesting objects in number theory, as well as in the Fourier analysis on finite fields. These sums are also obviously related to the *Jacobi theta function*, defined for complex $z \in \mathbb{C}$ and $\tau \in \mathbb{H} := \{x + iy \in \mathbb{C} : y > 0\}$ (the Poincaré upper half-space, a model for the hyperbolic plane) by the series

$$\theta(z, \tau) := \sum_{n=0}^{\infty} e^{\pi i \tau n^2 + 2\pi i z n}$$

If you plot the sums for a large number of values of N , given an irrational α or a rational with large denominator, you see “curlicues” as in the following pictures:



Theta sums with $\alpha \simeq \pi$, and $N = 10000$, 100000 and 1000000 .

Observe the axis: the sums are of the order of \sqrt{N} , as typical trajectories of a random walk!

ex: You may also explore what happens with other exponents, such as \sqrt{n} , and get interesting patterns or phenomena.

6 Simple orbits and perturbations

6.1 Topological fixed point theorems

To find periodic points of a transformation $f : X \rightarrow X$, namely fixed points of its iterates f^n , may be difficult. For example, when $f(x)$ is a polynomial of degree $d > 1$, its iterates are polynomials of exponentially growing degree.

Fixed point theorems in intervals. In real dimension one, connected and convex sets coincide, and are called intervals. This “miracle” is responsible for two very simple criteria to prove the existence of fixed points of continuous interval transformations. They say that if a compact interval is squeezed or stretched, at least one of its points remains fixed.

Theorem 6.1 (fixed point theorem for intervals). *Let $f : I \rightarrow \mathbb{R}$ be a continuous transformation defined in an interval $I \subset \mathbb{R}$.*

- i) If $J \subset I$ is a compact interval such that $f(J) \subset J$, then f has a fixed point in J .*
- ii) If $J \subset I$ is a compact interval such that $J \subset f(J)$, then f has a fixed point in J .*

The proof is an elementary application of Bolzano theorem to the continuous function $f(x) - x$. A more abstract proof, which can be generalized to higher dimension (with the help of some non-trivial algebraic topology), is as follows. Suppose that f has no fixed points in J . Then the function

$$g(x) = \frac{f(x) - x}{|f(x) - x|}$$

(which makes sense if the denominator does not vanish) would define a continuous map of an interval (J itself in case i) or some sub-interval of J in case ii), which is a connected space, onto the disconnected space $\{-1, 1\}$.

ex: Prove theorem 6.1.

ex: Find examples of continuous functions $f : I \rightarrow I$ and non-compact intervals J such that $f(J) \subset J$ or $J \subset f(J)$ which do not contain fixed points of f .

Other topological fixed point theorems. In higher, but finite, dimension, part i) of theorem 6.1 generalizes as

Theorem 6.2 (Brouwer). *A continuous map $f : D \rightarrow D$ of the closed unit disk $D \subset \mathbb{R}^{N+1}$ into itself has a fixed point.*

The idea is that if a continuous map $f : D \rightarrow D$ had no fixed point, the same formula as above (associating to each point x of the disk the intersection between the ray passing through x and $f(x)$ and the boundary sphere) would define a continuous map $g : D \rightarrow \partial D$ from the disk onto the unit N -sphere $\partial D = S^N$. That such a map cannot exist is quite clear intuitively, but needs some non-trivial algebraic topology to rigorously prove it (the N -th homotopy group of the $(N+1)$ -disk is trivial, while that of the N -sphere is not!).

In infinite dimension, one has the

Theorem 6.3 (Schauder-Tychonov). *A continuous transformation $f : K \rightarrow K$ of a compact and convex subset $K \subset X$ of a Banach space X (or of a locally convex topological vector space) has a fixed point.*

6.2 Dynamics of contractions

The simplest dynamical systems are contractions.

Contractions. Let (X, d) be a metric space. A map $f : X \rightarrow X$ is called *contraction* (or λ -contraction if one wants to keep track of the constant λ) if it is Lipschitz and has Lipschitz constant $\lambda < 1$, i.e. if there exists a $0 \leq \lambda < 1$ such that for all $x, x' \in X$

$$d(f(x), f(x')) \leq \lambda \cdot d(x, x') \quad (6.1)$$

Clearly, a constant transformation, sending any $x \in X$ into $f(x) = p$, is trivially a contraction. A linear homogeneous transformation $f(z) = \lambda z$ of the complex plane \mathbb{C} or of the real line \mathbb{R} is a contraction provided $|\lambda| < 1$. Observe that a contraction, as any Lipschitz map, is continuous (take $\delta = \varepsilon/\lambda$ in the ε - δ definition).

e.g. Smooth contractions. By the mean value theorem, a continuously differentiable transformation $f : \mathbb{R}^N \rightarrow \mathbb{R}^N$ (or of a convex subset $X \subset \mathbb{R}^N$) is a contraction provided there exists a positive $\lambda < 1$ such that $|f'(x)| \leq \lambda$ for any $x \in \mathbb{R}^N$.

ex: Show that a contraction of a compact metric space X cannot be invertible, provided the space contains more than one point (compare the diameters X and $f(X)$)

ex: Give non-trivial (i.e. non constant) examples of contractions of

$$[0, 1] \quad [0, 1] \times [0, 1] \quad B_r(x) = \{y \in \mathbb{R}^N \text{ t.q. } d(x, y) < r\} \quad \mathbb{S}^1 = \{z \in \mathbb{C} \text{ t.q. } |z| = 1\}$$

Contraction principle. The dynamics of a contraction is described by the following fundamental theorem, which we state with all details.

Theorem 6.4 (Contraction principle/Banach fixed point theorem). *All trajectories of a contraction $f : X \rightarrow X$ of a metric space (X, d) are fundamental sequences, and the distance between any two trajectories tends to zero exponentially fast. If X is complete, then f admits one and only one fixed point p . The trajectory of any initial point $x_0 \in X$ converges exponentially fast to the fixed point, i.e.*

$$d(f^n(x), p) \leq C \lambda^n,$$

where $C > 0$ is a positive constant and $0 \leq \lambda < 1$ is the Lipschitz constant of f .

Proof. Let $f : X \rightarrow X$ be a λ -contraction. Let $x_0 \in X$ be any initial point, and let (x_n) be its trajectory, defined by the recursion $x_{n+1} = f(x_n)$. Iterating (6.1), one sees that

$$d(x_{k+1}, x_k) \leq d(x_1, x_0) \cdot \lambda^k$$

Using k times the triangular inequality and then the convergence of the geometric series of ratio $\lambda < 1$, we get

$$\begin{aligned} d(x_{n+k}, x_n) &\leq \sum_{j=0}^{k-1} d(x_{n+j+1}, x_{n+j}) \leq d(x_1, x_0) \cdot \sum_{j=0}^{k-1} \lambda^{n+j} \\ &\leq d(x_1, x_0) \cdot \lambda^n \cdot \sum_{j=0}^{\infty} \lambda^j \leq \frac{\lambda^n}{1-\lambda} \cdot d(x_1, x_0). \end{aligned}$$

This implies that (x_n) is fundamental, since we can make $\lambda^n \cdot d(x_1, x_0)/(1-\lambda)$ smaller than any $\varepsilon > 0$ choosing a sufficiently large $n = n(\varepsilon)$. Continuity of f implies that the limit $p = \lim_{n \rightarrow \infty} x_n$, which exists if X is complete, is a fixed point of f . Uniqueness is clear, for if p and p' were two different fixed points, then by (6.1) their distance $\delta = d(p, p') > 0$ would be $\leq \lambda\delta$, which is impossible if $\lambda < 1$. Again by (6.1) and finite induction, the distance between any two trajectories $x_n = f^n(x_0)$ and $x'_n = f^n(x'_0)$ decay as $d(x_n, x'_n) \leq \lambda^n \cdot d(x_0, y_0)$. In particular, the distance between an arbitrary trajectory and the fixed point p is bounded by $d(x_n, p) \leq \lambda^n \cdot d(x_0, p)$, proving our last assertion with $C = d(x_0, p)$. \square

ex: Show that a transformation $f : X \rightarrow X$ of a complete metric space (X, d) such that

$$d(f(x), f(x')) < d(x, x')$$

for all distinct $x, x' \in X$ may fail to have fixed points (think at a decreasing sequence forming a divergent series).

ex: Let $a > 0$ and $x_0 > 0$. Show that the sequence (x_n) defined by

$$x_{n+1} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right)$$

converges to \sqrt{a} . This is Babylonians-Heron method to approximate square roots (the sequence is a trajectory of the transformation $f(x) = (x + a/x)/2$, which is a contraction once restricted to the closed interval $[\sqrt{a}, \infty) = f(\mathbb{R}_+) \dots$)

Stability of contractions. A contraction $f : X \rightarrow X$ of a complete metric space X may be thought as a “machine” that computes the fixed point $p = \lim_{n \rightarrow \infty} f^n(x)$ starting with any initial guess $x \in X$.

We pose the question whether contractions are stable, in some sense to be specified. We want to decide whether a small perturbation of a contraction f , say $g : X \rightarrow X$, produces a fixed point p' near to p . The point is to decide what “small” means. If we only require something like $d_\infty(f, g) := \sup_{x \in X} d(f(x), g(x)) < \delta$, the transformation g needs not be a contraction, no matter how small δ is chosen (to see this, try to visualize a δ -neighbourhood of the graph of a contraction of an interval, and fit there the graph of a transformation g with arbitrarily large derivative). It is clear that we also need some control on the derivatives. One possibility is to assume that X has a linear and differentiable structure, e.g. is a subset of some Euclidean $X \subset \mathbb{R}^N$, and look for f and g smooth. The condition

$$\|f - g\|_{C^1} := \sup_{x \in X} \|f(x) - g(x)\| + \sup_{x \in X} \|f'(x) - g'(x)\| < \delta$$

clearly implies that, if f is a λ -contraction and $\delta < 1 - \lambda$, then also g is a contraction and has Lipschitz constant $\leq \lambda + \delta$. Simpler, however, is to formulate a stability result inside the class of contractions.

Theorem 6.5. *Let $f : X \rightarrow X$ be a λ -contraction of the complete metric space (X, d) , and let $p \in X$ be its fixed point. For every $\varepsilon > 0$ there exists some $0 < \delta < 1 - \lambda$ such that if $g : X \rightarrow X$ is a $(\lambda + \delta)$ -contraction at distance $d_\infty(f, g) < \delta$ from f , and if p' is the fixed point of g , then*

$$d(p, p') < \varepsilon$$

Proof. If p' is the fixed point of g , we know that $g^n(p) \rightarrow p'$ when $n \rightarrow \infty$. By triangle inequality we see that

$$\begin{aligned} d(p, p') &\leq \sum_{n=0}^{\infty} d(g^{n+1}(p), g^n(p)) \leq d(g(p), p) \cdot \sum_{n=0}^{\infty} (\lambda + \delta)^n \\ &\leq \delta \cdot \sum_{n=0}^{\infty} (\lambda + \delta)^n = \frac{\delta}{1 - (\lambda + \delta)} \end{aligned}$$

and this quantity is $< \varepsilon$ provided that δ is sufficiently small. □

Equivalence classes of linear contractions. Linear contraction of the real line, transformations $f(x) = \lambda x$ with $|\lambda| < 1$, also provide a simple example of how to use the dynamics to build a conjugation between two transformations.

Theorem 6.6. *Let $f : x \mapsto \alpha x$ and $g : x \mapsto \beta x$ be two linear non-trivial contractions of \mathbb{R} . They are topologically conjugated iff α and β share the same sign.*

Proof. Assume first that $0 < \alpha < 1$ and $0 < \beta < 1$. The origin is the common fixed point. The set $A = [-1, -\alpha) \cup (\alpha, 1]$ is a “fundamental domain” for the action of f on the punctured real line $\mathbb{R}^\times := \mathbb{R} \setminus \{0\}$, in the sense that for any $x \in \mathbb{R} \setminus \{0\}$ there exists a unique time $n(x) \in \mathbb{Z}$ such that $f^{n(x)}(x) \in A$. Similarly, a fundamental domain for the action of g on \mathbb{R}^\times is $B = [-1, -\beta) \cup (\beta, 1]$. Let $H : A \rightarrow B$ be any homeomorphism such that $H(-1) = -1$, $H(-\alpha) = -\beta$, $H(\alpha) = \beta$ and $H(1) = 1$ (for example, an affine homeomorphism). It is easy to check that the recipe

$$h(x) = \begin{cases} 0 & \text{se } x = 0 \\ g^{-n(x)}(H(f^{n(x)}(x))) & \text{if } x \neq 0 \end{cases}$$

defines a homeomorphism $h : \mathbb{R} \rightarrow \mathbb{R}$. Since $n(x) = n(f(x)) + 1$ (why?), we see that

$$\begin{aligned} (h \circ f)(x) &= g^{-n(f(x))} \left(H \left(f^{n(f(x))} f(x) \right) \right) = g^{-n(x)+1} \left(H \left(f^{n(x)-1} (f(x)) \right) \right) \\ &= g \left(g^{-n(x)} \left(H \left(f^{n(x)}(x) \right) \right) \right) = (g \circ h)(x) \end{aligned}$$

and therefore h is a topological conjugation between f and g .

The case when $-1 < \alpha, \beta < 0$ is analogous. On the other side, it is not difficult to see that the contractions $x \mapsto \alpha x$ and $x \mapsto -\alpha x$, with $\alpha \neq 0$, having opposed orientations, cannot be conjugated. \square

The result is that non-trivial linear contractions of the real line fit into two classes of topological conjugated transformations, those that preserve the orientation and those that invert the orientation.

It is important to observe that a conjugation h between two contractions $f : x \mapsto \alpha x$ and $g : x \mapsto \beta x$ cannot be differentiable, unless $\alpha = \beta$. Indeed, if $f = h^{-1} \circ g \circ h$ and if h is differentiable, then the chain rule implies that $f'(0) = g'(0)$, hence that $\alpha = \beta$.

ex: Show that the linear contractions $x \mapsto \alpha x$ and $x \mapsto -\alpha x$, with $\alpha \neq 0$, cannot be conjugated (observe that a conjugation is a homeomorphism of the line, in particular monotone).

6.3 Linear maps

Linear systems are the only dynamical systems we can explicitly solve. They serve as models of the local behavior of generic systems near a fixed point.

Linear maps. A linear transformation of the Euclidean vector space \mathbb{R}^N (which we may think equipped with the standard Euclidean structure) is an endomorphism $f : \mathbb{R}^N \rightarrow \mathbb{R}^N$, defined, in the canonical basis, by a square matrix $A = (a_{ij}) \in \text{Mat}_{n \times n}(\mathbb{R})$, according to

$$f(x) = Ax$$

Here, we think at $x = (x_1, x_2, \dots, x_N)^\top$ as a column vector, and therefore at Ax as the usual row-by-columns product between matrices. Thus, the map is defined, in coordinates, by $x_i \mapsto \sum_j a_{ij}x_j$.

A linear change of coordinates $h : x \mapsto y = Ux$, with $U \in GL_N(\mathbb{R})$ an invertible square matrix, defines a linear (and therefore topological) conjugation between f and the linear transformation $g : y \mapsto By$, defined by the matrix $B = UAU^{-1}$. Thus, we are free to change coordinates.

Observe that the origin is a fixed point of f , i.e. $f(0) = 0$, and it is the unique fixed point iff the matrix $A - I$ is invertible, i.e. if 1 is not an eigenvalue of A .

If x is an eigenvector of A with eigenvalue λ , i.e. a non-trivial solution of the homogeneous equation $Ax = \lambda x$, then iterations are simple. Indeed, $f^n(x) = \lambda^n x$, and therefore the asymptotic behaviour of the trajectory of the eigenvector x depends on the absolute value of its eigenvalue. Trajectories converge to the origin when $|\lambda| < 1$, and diverge when $|\lambda| > 1$. In the exceptional cases with $\lambda = \pm 1$, we have a fixed point or a periodic point with period 2.

In order to understand the possible global pictures, we start with the smallest non-trivial case.

Linear maps in the plane. The simplest non-trivial case is that of a linear map endomorphism of the plane \mathbb{R}^2 , defined, in the canonical basis, by a two-by-two real square matrix A . The qualitative behaviour of trajectories depends on the eigenvalues of A , and on their geometric multiplicity. Remember that the eigenvalues λ_{\pm} are the roots of the characteristic polynomial

$$\det(\lambda I - A) = \lambda^2 - (\text{tr} A)\lambda + \det A,$$

which is a degree two polynomial with real coefficients, and therefore they are a couple of real numbers $\lambda_{\pm} = \alpha \pm \beta$, possibly overlapping, or a couple of complex conjugate numbers $\lambda_{\pm} = \alpha \pm i\omega$.

If the matrix A is diagonalizable (as a real matrix in a real vector space), i.e. admits two eigenvalues λ_{\pm} (possibly equal) and two linearly independent eigenvectors v_{\pm} , then the system is linearly conjugated to a diagonal system

$$f(x, y) = (\lambda_+ x, \lambda_- y).$$

If both eigenvalues have absolute values $|\lambda_{\pm}| < 1$, then f is a contraction and the orbit of any point converges (exponentially fast) to $f^n(x, y) \rightarrow 0$. Trajectories move along curves

$$y = Cx^{\alpha},$$

for some constant C and some exponent $\alpha = \log |\lambda_-| / \log |\lambda_+| > 0$. The basin of attraction of the origin is the whole plane, and the origin, which is an attractng fixed point, is called a *stable node*, or *sink*.

If both eigenvalues have absolute values $|\lambda_{\pm}| > 1$, then the inverse f^{-1} is a contraction, all backward trajectories converge to $f^{-n}(x, y) \rightarrow 0$ as $n \rightarrow \infty$, and all forward tarjectories of points different from the origin diverge, i.e. $|f^n(x, y)| \rightarrow \infty$ as $n \rightarrow \infty$. The basin of attraction of the origin is $\{0\}$ itself, and the origin is called a *unstable node*, or *source*.

If one of eigenvalue has $|\lambda_-| < 1$ and the other $|\lambda_+| > 1$, what happens is the following: trajectories starting at the “stable line” $E^- = \mathbb{R}v_- \approx \{(0, y)\}$, the eigenspace of the eigenvalue λ_- , converge to the origin, while trajectories starting at the “unstable line” $E^+ = \mathbb{R}v_+ \approx \{(x, 0)\}$, the eigenspace of the eigenvalue λ_+ , diverge. A generic trajectory, starting at a point which does not belong to $E^- \cup E^+$, i.e. (x, y) with both $x \neq 0 \neq y$, also diverge (since the y coordinate decays but the x coordinate explodes), moving along curves

$$y = Cx^{\beta},$$

for some constant C and some exponent $\beta = \log |\lambda_-| / \log |\lambda_+| < 0$. The origin is then called a *saddle*, and the linear map *hyperbolic*.

The next case is when A has only one eigenvalue λ , with geometric multiplicity one (i.e. admits just a one-dimensional family of eigenvectors). It can be shown that the system is linearly conjugated with

$$f(x, y) = (\lambda x + y, \lambda y).$$

(any eigenvector is proportional to $(1, 0)$). One easily check that iterations of this map are

$$f^n(x, y) = \lambda^{n-1}(\lambda x + ny, \lambda y).$$

Therefore, if $|\lambda| < 1$ all trajectories converge to the origin, which is then called a *degenerate stable node*. If $|\lambda| > 1$, then all trajectories of points different from the origin diverge. The origin is then called *degenerate source*. It is clear, however, that this is not a generic situation. A small (generic) perturbation of the matrix leads to one of the previous cases, or to the following case.

Finally, it may happen that the characteristic polynomial has no real roots, but a couple of complex conjugated roots $\lambda_{\pm} = \rho e^{\pm i\theta}$, for some $\rho = |\lambda_{\pm}| > 0$ and $\theta \notin \pi\mathbb{Z}$. This means that the

complexification of A , the linear operator $A^{\mathbb{C}} : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ defined in the canonical basis by the same matrix as A , admits two linearly independent eigenvectors v_{\pm} , corresponding to the two complex eigenvalues λ_{\pm} . Moreover, since $A = \bar{A}$, we may take $v_- = \overline{v_+}$. But then, in the basis of $\mathbb{R}^2 \subset \mathbb{C}^2$ defined by $e_1 = (v_+ + v_-)/2$ and $e_2 = (v_+ - v_-)/2i$, the map $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is induced/defined by the matrix

$$B = \rho \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}.$$

This means that f is a (counter-clockwise) rotation by an angle θ followed by a homothety/scaling with ratio ρ . Iterations of B are simply

$$B^n = \rho^n \begin{pmatrix} \cos(n\theta) & -\sin(n\theta) \\ \sin(n\theta) & \cos(n\theta) \end{pmatrix}$$

To understand trajectories, it is easier to identify the plane with the complex line $\mathbb{R}^2 \approx \mathbb{C}$, and use polar coordinates $(x, y) \approx x + iy = re^{i\varphi}$. Then $f^n(re^{i\varphi}) = r\rho^n e^{i(\varphi+n\theta)}$, and therefore trajectories move along logarithmic spirals

$$r = Ce^{\gamma\varphi},$$

for some constant C and some exponent $\gamma = (\log \rho)/\theta$, which may be positive or negative, or along circles $r = C$, if it happens that $\rho = 1$. In particular, if $|\lambda_{\pm}| = \rho < 1$, then all trajectories converge to the origin $f^n(x, y) \rightarrow 0$ as $n \rightarrow \infty$, which is then called a *stable focus*. If $|\lambda_{\pm}| = \rho > 1$, then the trajectories of all points different from the origin diverge, and the origin is then called an *unstable focus*.

ex: Describe what happens in the exceptional situation when $f(x, y) = (x + y, y)$, i.e. the only eigenvalue is 1 and it has geometric multiplicity one.

General linear maps, Jordan normal form. Let $f(x) = Ax$ be a linear system defined by a real $n \times n$ matrix A , and consider its complexification, i.e. the linear operator $A : \mathbb{C}^N \rightarrow \mathbb{C}^N$, acting on $\mathbb{C}^N = \mathbb{R}^N \oplus i\mathbb{R}^N$ according to $A(x + iy) := Ax + iAy$. According to the Jordan normal form theorem, the complexified linear space is a direct sum $\mathbb{C}^N = \bigoplus E_{\lambda}$ of *generalized eigenspaces*, or *root spaces*, E_{λ} , which are invariant under A and where the action of A is

$$\lambda I + N$$

where λ is the eigenvalue, and N is a nilpotent operator. More precisely, the matrix which represents the restriction of A on $E_{\lambda} \subset \mathbb{C}^N$ is a Jordan block

$$J = \begin{pmatrix} \lambda & 1 & & & \\ & \lambda & 1 & & \\ & & \ddots & \ddots & \\ & & & \lambda & 1 \\ & & & & \lambda \end{pmatrix} \quad (6.2)$$

(empty entries are all zero). Moreover, generalized eigenspaces with non-real eigenvalues come in pairs of generalized eigenspaces E_{λ} and $E_{\bar{\lambda}}$, whose vectors are related by a complex conjugation. As in the two-dimensional situation, one can then construct an invariant subspace $E_{\lambda, \bar{\lambda}} \subset \mathbb{R}^N$ where the action of A is given by the Jordan block

$$\begin{pmatrix} \rho R_{\theta} & I & & & \\ & \rho R_{\theta} & I & & \\ & & \ddots & \ddots & \\ & & & \rho R_{\theta} & I \\ & & & & \rho R_{\theta} \end{pmatrix} \quad (6.3)$$

where

$$\rho R_{\theta} = \rho \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} = \begin{pmatrix} \alpha & -\beta \\ \beta & \alpha \end{pmatrix} \quad \text{and} \quad I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

and $\lambda = \alpha + i\beta = \rho e^{i\theta}$.

The conclusion is that the phase space \mathbb{R}^N of the real linear system $f(x) = Ax$ splits as a direct sum of invariant subspaces, E_λ or $E_{\lambda, \bar{\lambda}}$, where A acts as (6.2) or as (6.3), respectively.

It is clear that the asymptotic behavior of iterations of A on each root space depends on the absolute value $|\lambda|$ of the corresponding eigenvalue. Indeed, one can write the total space as a direct sum of three invariant subspaces

$$\mathbb{R}^N = E^- \oplus E^0 \oplus E^+$$

where the *stable space* E^- is the direct sum of those root spaces with $|\lambda| < 1$, the *unstable space* E^+ is the direct sum of those root spaces with $|\lambda| > 1$, and finally the *neutral space* E^0 is the direct sum of those root spaces with $|\lambda| = 1$.

One can then show that the restriction of f to E^- is eventually contracting (i.e. some power is contracting), and therefore $f^n(x) \rightarrow 0$ as $n \rightarrow \infty$ if $x \in E^-$. This happens because the exponential contraction dominates the nilpotent part of each Jordan block for sufficiently large n , and therefore $\|A^n v\| \leq C\mu^n \|v\|$ for some constants C and $\mu < 1$ and every $v \in E^-$. Similarly, one shows that the inverse of the restriction of f to E^+ is eventually contracting.

A linear map is called *hyperbolic* if all the eigenvalues have $|\lambda| \neq 1$, i.e. if the spectrum of the complexification is disjoint from the unit circle of the complex plane. This means that the phase space splits as a direct sum

$$\mathbb{R}^N = E^- \oplus E^+$$

of a stable and an unstable subspace. If an hyperbolic map has eigenvalues with absolute values both $|\lambda| < 1$ and $|\lambda| > 1$, then the origin is called a *saddle*.

6.4 Order of the line and trajectories

The order of the real line cause restrictions of possible trajectories of monotone transformations.

Increasing maps of the interval. Let $f : I \rightarrow I$ be a continuous increasing map of the interval $I \subset \mathbb{R}$. Any trajectory $(x_n)_{n \in \mathbb{N}_0}$ is monotone, increasing or decreasing (depending whether $f(x) > x$ or $f(x) < x$, respectively). The monotone trajectory may converge, i.e. $x_n \rightarrow p$ to some fixed point, if bounded, or may diverge $x_n \rightarrow \pm\infty$, is unbounded. In particular, if the interval I is compact, the second possibility is excluded, and all trajectories converge to some fixed point. In this case, there exists a not-empty compact set $F \subset I$ made of fixed points, and any x in each connected component A_k of $I \setminus F = \bigcup_k A_k$ has a trajectory contained in A_k which converge to some point $x_\infty \in \partial A_k$.

ex: Show that a homeomorphism $f : I \rightarrow I$ of an interval $I \subset \mathbb{R}$ cannot have periodic points of period larger than 2. When does it have periodic points of period 2?

ex: Let $I \subset \mathbb{R}$ be compact interval and $f : I \rightarrow I$ a continuous and increasing function. Show that any trajectory converges to a fixed point. Discuss the dynamics of f .

ex: Discuss the dynamics of a continuous and decreasing map $f : I \rightarrow I$ of a compact interval $I \subset \mathbb{R}$ (observe that if f is decreasing then f^2 is increasing).

ex: (difficult!) Let $I \subset \mathbb{R}$ be an interval and $f : I \rightarrow I$ and $g : I \rightarrow I$ be two homeomorphisms of I without fixed points. Show that they are topologically conjugated.

Sharkovskii order. The order of the real line also implies restrictions on the possible periods of a map. A striking result by Alexander N. Sharkovskii³⁰ says that there exists an order \prec on the naturals, which looks like

$$\begin{aligned} 1 &\prec 2 \prec 2^2 \prec 2^3 \prec \dots \prec 2^m \prec \dots \prec 2^k \cdot (2n-1) \prec \dots \\ \dots &\prec 2^k \cdot 3 \prec \dots \prec 2 \cdot 3 \prec \dots \prec 2n-1 \prec \dots \prec 9 \prec 7 \prec 5 \prec 3 \end{aligned}$$

³⁰A.N. Sharkovskii, Co-existence of cycles of a continuous mapping of the line into itself, *Ukrainian Math. J.* **16** (1964), 61-71.

such that if a continuous transformation $f : \mathbb{R} \rightarrow \mathbb{R}$ has an orbit of period k and if $j \prec k$ then it also has an orbit of period j . In particular, the existence of an orbit of period 3 implies the existence of orbits of all periods!

ex: Try to figure a transformation of the real line with an orbit of period 3.

6.5 Local analysis: attracting and repelling fixed points

Differentiability of a transformation and the contraction principle helps to understand the trajectories of points which are near to the fixed or periodic points.

Attracting and repelling fixed points. Let $f : X \rightarrow X$ be a transformation of class \mathcal{C}^1 defined in some open subset $X \subset \mathbb{R}^N$, and let $p \in X$ be a fixed point of f .

We say that the fixed point p is *attracting* if its basin of attraction $W^s(p)$ is a neighbourhood of p , i.e. if p admits a neighbourhood B such that $f^n(x) \rightarrow p$ for all $x \in B$. The following criterium is a simple consequence of the contraction principle.

Theorem 6.7. *If $|f'(p)| < 1$, then p is attracting.*

Proof. By continuity of f' , there exist $\lambda < 1$ and a ball $B = B_\varepsilon(p)$ around p such that $|f'(x)| < \lambda$ for all $x \in B$. By the mean value theorem, $f(\overline{B}) \subset \overline{B}$, since if $d(x, p) \leq \varepsilon$ then

$$d(f(x), p) \leq \lambda \cdot d(x, p) < \varepsilon$$

Moreover, the mean value theorem also implies that $d(f(x), f(x')) \leq \lambda \cdot d(x, x')$ if $x, x' \in \overline{B}$. Thus, $f|_{\overline{B}} : \overline{B} \rightarrow \overline{B}$ is a contraction, and the contraction principle says that trajectories of all points $x \in \overline{B}$ converge exponentially to p . \square

We say the fixed point p is *repelling* if it admits a neighbourhood B such that the trajectory of any $x \in B$, different from p , leaves B in finite time, i.e. $f^n(x) \notin B$ for some $n \geq 1$. The following criterium use the order of the real line.

Theorem 6.8. *Let $f : X \rightarrow X$ be a transformation of class \mathcal{C}^1 defined in some open interval $X \subset \mathbb{R}$. If $|f'(p)| > 1$, then p is repelling.*

Proof. By continuity of f' , there exist $\lambda > 1$ and an interval $B = [p - \varepsilon, p + \varepsilon]$ around p such that $|f'(x)| > \lambda$ for all $x \in B$. Also observe that f is strictly increasing or decreasing, depending on the sign of $f'(p)$, and therefore sends bijectively intervals onto intervals. take a point $x \in B$ different from p , and suppose that $f^k(x) \in B$ for all times $0 \leq k \leq n$. The chain rule implies that the derivative of f^n at points c between p and x grow exponentially, since

$$|(f^n)'(c)| = |f'(f^{n-1}(c))| \cdot |f'(f^{n-2}(c))| \cdots |f'(c)| > \lambda^n$$

The mean value theorem implies that n cannot be arbitrarily large, since

$$d(p, f^n(x)) \geq \lambda^n \cdot d(p, x) \quad \text{and} \quad d(p, f^n(x)) \leq \varepsilon$$

are not compatible for large n . Thus, there exists a time $n \geq 1$ such that $f^n(x) \notin B$. \square

It must be said that this result is local, it does not say anything about the basin of attraction of p . Also, the condition $|f'(p)| > 1$ is not sufficient to establish a similar result in higher dimension (there may be directions where f dilates distances, and directions where it contracts distances ...)

ex: Show by examples that the basin of attraction of a repelling fixed point p can be larger than $\{p\}$.

ex: Find a good definition of *attracting periodic orbit* (observe that the derivative of f^n is constant along a periodic orbit of period n , then consider iterations of $f^n \dots$)

ex: Consider the family of quadratic maps

$$x \mapsto \lambda x^2$$

depending on the parameter λ . Find the basin of attraction of the fixed point $p = 0$, and describe the speed of convergence of convergent trajectories.

ex: If p is a fixed point of $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $f'(p) = 1$, then everything can happen! The basin of attraction of p can be a neighbourhood of p , or just $\{x\}$, or may contain an half-neighbourhood like $[p, p + \varepsilon) \dots$

Consider the examples

$$x \mapsto x \pm x^3 \quad \text{e} \quad x \mapsto x \pm x^2$$

and find others.

The quadratic family. The *quadratic family* is the family of transformations of the unit interval $f_\lambda : [0, 1] \rightarrow [0, 1]$, defined according to

$$f_\lambda(x) = \lambda x(1 - x)$$

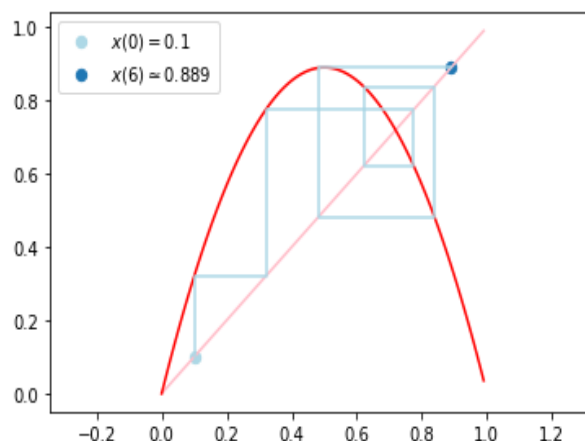
Here the parameter λ takes values in the interval $[0, 4]$. It is also called *logistic* (from the French “logement”), since it is a model of population growth in a limited environment, x being the relative population, the quotient of the actual population over the maximal allowed population.

Fixed points are 0, which is attracting for $0 \leq \lambda < 1$, and $p_\lambda = \frac{\lambda-1}{\lambda}$, which appears when $\lambda > 1$ (remember that our phase space is only the unit interval and not the whole real line) and is attracting when $1 < \lambda < 3$.

If $\lambda \in [0, 1]$ then all trajectories converge to 0. Indeed, trajectories are bounded and decreasing sequences, and 0 is the unique fixed point.

If $\lambda \in (1, 3]$ then all trajectories converge to p_λ . This is not so obvious.

What is really interesting is to observe what happens for increasing values of $\lambda > 3$. You may take a look at my applet in <http://w3.math.uminho.pt/~scosentino/salbestiario.html>.



Cobweb plot of the logistic map, for $\lambda \simeq 3.56$.

Convergence for Newton method. Let $F \in \mathbb{R}[x]$ be a polynomial with real coefficients. Newton method to find the roots of F , i.e. to solve the equation $F(x) = 0$, consists in choosing a first approximation x_0 , and then iterate

$$x_{n+1} = x_n - \frac{F(x_n)}{F'(x_n)}.$$

This means that we try to refine our bet x_n using the linear approximation (first order Taylor)

$$F(x) \simeq F(x_n) + F'(x_n) \cdot (x - x_n)$$

Clearly, we may iterate provided the derivative stays away from zero in a neighbourhood of the root we want to approximate.

It is clear that if the sequence (x_n) converges to some p , and if $F'(p) \neq 0$, then the limit $p = \lim_{n \rightarrow \infty} x_n$ is a root of the polynomial F . Conversely, if p is a root of F , and if $F'(p) \neq 0$ (so that it is also different from zero in a neighbourhood of p), then p is a fixed point of the map

$$x \mapsto f(x) := x - \frac{F(x)}{F'(x)}$$

The derivative of f at p is

$$f'(p) = 1 - \frac{(F'(p))^2 - F(p) F''(p)}{(F'(p))^2} = 0$$

Therefore, p is an attracting fixed point of f : the trajectory of any initial guess x_0 sufficiently close to p converges to p .

Indeed, since the derivative is $f'(p) = 0$, any root of F is a super-attracting fixed point of f , and the convergence is much better than exponential.

Theorem 6.9. *Let p be a non-critical root of the polynomial $F \in \mathbb{R}[x]$, i.e. a root where $F'(p) \neq 0$. Then Newton's iterations starting from any x_0 sufficiently near the root p converge to this root, and the convergence is "quadratic", i.e. the error $\varepsilon_n = |x_n - p|$ decreases as*

$$\varepsilon_{n+1} \leq K \cdot \varepsilon_n^2$$

for some $K > 0$.

Proof. We may assume, without loss of generality, that the root we are looking for is the origin, so that $F(0) = 0$. Now, suppose we are at x_n after n iterations. Taylor's formula with Lagrange estimate of the error around x_n says that

$$F(x) = F(x_n) + F'(x_n) \cdot (x - x_n) + \frac{1}{2} F''(y) \cdot (x - x_n)^2$$

for some y between x and x_n . Taking $x = 0$ (the root!) and dividing by $F'(x_n)$ we get

$$0 = F(0) = F(x_n) - F'(x_n) \cdot x_n + \frac{1}{2} F''(y) \cdot x_n^2$$

and therefore

$$x_n - \frac{F(x_n)}{F'(x_n)} = \frac{1}{2} \frac{F''(y)}{F'(x_n)} x_n^2$$

But the l.h.s. is x_{n+1} , so that

$$x_{n+1} = \frac{1}{2} \frac{F''(y)}{F'(x_n)} x_n^2$$

Since $F'(0) \neq 0$ (and polynomials have continuous derivatives), there is an interval $I =]-\varepsilon, \varepsilon[$ around the root 0 where $M = \sup_{x \in I} |F''(x)| < \infty$ and $\delta = \inf_{x \in I} |F'(x)| > 0$. Let $K = M/2\delta$. There follows that the distance $\varepsilon_n = |x_n - 0|$ between the n -th iterate and the root satisfies the iterative bound

$$|\varepsilon_{n+1}| \leq K \cdot |\varepsilon_n|^2$$

□

ex: Check that Newton's method applied to the quadratic polynomial $z^2 - a$, with $a > 0$, corresponds to Heron's algorithm.

ex: Estimate $\sqrt{17}$.

ex: Write down Newton's repice to solve $z^n - a = 0$, with $a > 0$ and $n \geq 2$.

ex: Use Newton's method to estimate the roots of

$$z^2 + 1 + z \quad z^3 - z - 1 \quad z^5 + z + 1 \quad z^3 - 2z - 5$$

Linearization in the complex plane. Let $f : \mathbb{C} \rightarrow \mathbb{C}$ be a rational function defined in the Riemann sphere $\mathbb{C} = \mathbb{C} \cup \{\infty\}$. Any fixed point p has its basin of attraction B_p . Looking for fast methods to compute the iterates, in 1871 E. Schröder had the idea to look for local conformal conjugations of f with simpler rational functions, like affine functions $g : z \mapsto \lambda z$. The method amounts to solve the functional equation

$$h \circ f|_{B_p} = g \circ h,$$

where $h : B_p \rightarrow B$ is an holomorphic function. E. Schröder, G. Koenig and J.H. Poincaré solved the problem with $|\lambda| \neq 1$, and then Carl S. Siegel solved the case $|\lambda| = 1$ around 1940.

Theorem 6.10 (Koenigs). *Let z_0 be a fixed point of f with multiplier $f'(z_0) = \lambda$ such that $|\lambda| \neq 0, 1$. Then there exists a conformal map ϕ , unique up to a non-zero factor, from a neighbourhood of z_0 onto a neighbourhood of 0 such that $\phi \circ f = \lambda \cdot \phi$.*

Proof. We assume that z_0 is attracting, i.e. $|\lambda| < 1$, since the repelling case follows considering the local inverse of f . Also, after conjugation, we can assume that $z_0 = 0$, hence the map has the form

$$f(z) = \lambda z + a_2 z^2 + \dots$$

Now define $\phi_n(z) = f^n(z)/\lambda^n$. There exists a $\delta > 0$ and a constant $c < |\lambda| < 1$ such that, for $|z| < \delta$,

$$|\phi_{n+1}(z) - \phi_n(z)| \leq k \cdot (c/|\lambda|)^n$$

for some $k > 0$. Hence the sequence of holomorphic functions ϕ_n converges uniformly in a small ball around 0. The functional equation $\phi \circ f = \lambda \cdot \phi$ follows immediately from its definition.

Comparing coefficients it is easy to see that any conjugation of $z \mapsto \lambda z$ to itself is a constant multiple of the identity, as long as $|\lambda| \neq 0, 1$. Uniqueness follows. \square

Theorem 6.11 (Böttcher). *Let z_0 be a superattracting fixed point of f , where*

$$f(z) = z_0 + a_p(z - z_0)^p + \dots$$

with $a_p \neq 0$ and $p \geq 2$. Then there exists a conformal map ϕ , unique up to multiplication by a $(p-1)$ -root of unity, from a neighbourhood of z_0 onto a neighbourhood of 0 such that $\phi \circ f = \phi^p$.

Proof. (sketch) We can assume that $z_0 = 0$ and that $a_p = 1$. As in Koenigs proof, we look for the conjugation as a limit of the functions

$$\phi_n(z) = f^n(z)^{p^{-n}}.$$

It can be shown that the ϕ_n converge uniformly in some sufficiently small ball around 0, and the functional equation follows from the definition. Uniqueness, up to a $(p-1)$ -root of unity, can be checked comparing power series. \square

6.6 Transversality and bifurcations

Transversality. Let $f : I \rightarrow \mathbb{R}$ be a transformation of class \mathcal{C}^1 defined in some interval $I \subset \mathbb{R}$, and let p be a fixed point of f . If $f'(p) \neq 1$, then this fixed point is “isolated”, i.e. it is the unique fixed point of f in some neighbourhood of p . Indeed, a fixed point is a solution of

$$F(x) = f(x) - x = 0$$

Now, if $f'(p) \neq 1$ then $F'(p) \neq 0$. The inverse function theorem then says that F is invertible in a neighbourhood B of p , and this implies that p is the unique zero of F in B , so that p is the unique fixed point of f in B .

Fixed points satisfying the condition $f'(p) \neq 1$ are called *transversal*, because the tangent to the graph of f at p is transversal to the (tangent to the) graph of the identity function.

Persistence. The condition $f'(p) \neq 1$ is an open condition, and this suggests that it may stable under small perturbations of f .

Theorem 6.12. Let $f : I \rightarrow \mathbb{R}$ be a transformation of class \mathcal{C}^1 , and p be a transversal fixed point of f . Then all transformations $g : I \rightarrow \mathbb{R}$ sufficiently \mathcal{C}^1 -near to f have one, and only one, fixed point in some neighbourhood of p , which is also transversal.

Proof. Let $g = f - h$ be a perturbation of f , with $\|h\|_{\mathcal{C}^1} = \|h\|_{\infty} + \|h'\|_{\infty} < \delta$. A fixed point of g is a solution of $g(x) - x = 0$, i.e. of

$$F(x) = h(x)$$

if we define $F(x) = f(x) - x$. We know that F in some neighbourhood B' of p , hence a fixed point of g inside B is a solution of $x = (F^{-1} \circ h)(x)$, which means a fixed point of $F^{-1} \circ h$. The strategy, now, it to show that $F^{-1} \circ h$ is a contraction in some neighbourhood of p . If the closed neighborhood $B = \overline{B_r(p)}$ is sufficiently small, then the inverse of F has bounded derivative, say $|(F^{-1})'(x)| < M$ in $F(B)$. If δ is sufficiently small, then the derivative $|(F^{-1} \circ h)'(x)| < M \cdot \delta$ is uniformly $\leq \lambda := M\delta < 1$ in B , and therefore $F^{-1} \circ h$ has good chances to be a contraction. We must show that the image $(F^{-1} \circ h)(B)$ is contained in B . Now, given $x \in B$, triangular inequality, the mean value theorem and the chain rule, imply o

$$\begin{aligned} d((F^{-1} \circ h)(x), p) &\leq d(F^{-1}(h(x)), F^{-1}(h(p))) + d(F^{-1}(h(p)), p) \\ &\leq d(F^{-1}(h(x)), F^{-1}(h(p))) + d(F^{-1}(h(p)), F^{-1}(0)) \\ &\leq M \cdot \delta \cdot r + M \cdot \delta \end{aligned}$$

(where we used the fact that p is a fixed point of f) and this quantity is $< r$ whenever δ is sufficiently small. The contraction principle then says that a fixed point $p' \in B$ of g exists and is unique. The derivative of g at this point is δ -near the derivative of f in p , and therefore this fixed point is also transversal. \square

ex: Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a transformation of class \mathcal{C}^1 , and let p a periodic point of period n such that $(f^n)'(p) \neq 1$. Show that all transformations sufficiently \mathcal{C}^1 -near to f have a periodic point of period n near p . (consider the iterate f^n and apply the above theorem)

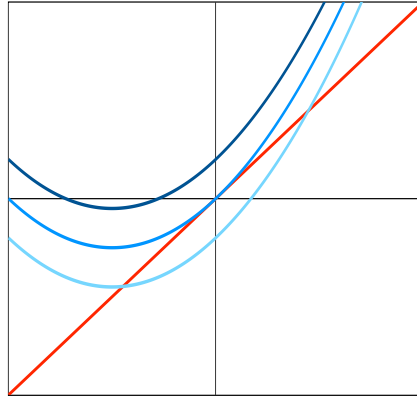
ex: Let $f : V \rightarrow \mathbb{R}^N$ be a transformation of class \mathcal{C}^1 defined in some open set $V \subset \mathbb{R}^N$, and let p be a fixed point of f . Transversality of p now means that the derivative (Jacobian) operator $f'(p)$ does not have 1 as an eigenvalue. State and prove the analogous of theorem 6.12 in this case.

Bifurcations. Non-transversal fixed points need not be persistent, and may disappear or change their nature under generic perturbations. This phenomenon is called *bifurcation*. The idea of bifurcation theory is to treat families f_λ of transformations, defined in some neighbourhood of a fixed or periodic point of $f = f_0$, and describe possible changes in the dynamics when the parameter λ varies.

Consider, for example, the family

$$f_\lambda(x) = x + x^2 - \lambda$$

defined in the real line. The origin is a non-transversal fixed point of f_0 , where $f'_0(0) = 1$. If $\lambda \neq 0$ is small, then f_λ has two fixed points $\pm\sqrt{\lambda}$, one repelling and the other attracting, if $\lambda > 0$, or none if $\lambda < 0$.



Graphs of $f(x) = x + x^2 - \lambda$, for $\lambda = -0.2, 0$ and 0.2 (different kinds of blue), compared with the diagonal (red).

The family

$$f_\lambda(x) = x + x^3 + \lambda x$$

shows a different behavior. The problem is to decide which phenomena are “generic”, and possibly “stable”, in some sense to be specified.

If we admit the existence of a sufficient number of derivatives, an arbitrary family of transformations with a non-transversal fixed point at the origin when $\lambda = 0$ is

$$\begin{aligned} f_\lambda(x) &= a_\lambda + b_\lambda x + c_\lambda x^2 + \dots \\ &= (a'\lambda + a''\lambda^2 + \dots) + (1 + b'\lambda + b''\lambda^2 + \dots)x + (c + c'\lambda + c''\lambda^2 + \dots)x^2 + \dots \end{aligned}$$

The generic case is when $c \neq 0$ (i.e. $f_0 = x + cx^2 + \dots$) and a generic perturbation has $a' \neq 0$ (i.e. the constant term of f_λ is different from zero when $\lambda \neq 0$). It is then not difficult to convince oneself that this family behaves qualitatively like the simpler family $f_\lambda(x) = x + x^2 - \lambda$ above. A small perturbation of f_0 may destroy the fixed point, in one direction, or create two new fixed points, in the other direction.

ex: State and prove the above result (observe that looking for roots of $f_\lambda(x) = x$, as function of λ , amount to to define look for functions $\lambda \mapsto x(\lambda)$ which satisfy $G(\lambda, x) = f_\lambda(x) - x = 0$, and this problem is treated by the implicit function theorem).

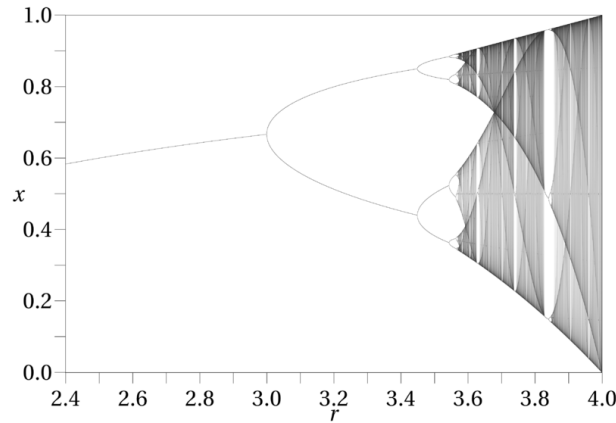
Period-doubling and Feigenbaum universality. Also interesting is the case of a family interval transformations f_λ such that f_0 has a fixed point at the origin with multiplies $f'_0(0) = -1$. Such a fixed point is transversal, hence persistent. Meanwhile, $(-1)^2 = 1$, and therefore the derivative of f_0^2 at 0 is $(f_0^2)'(0) = 1$. This says that the origin is not transversal as a fixed point of the second iterate f^2 . A small perturbation may produce periodic points of period 2 in a neighbourhood of the persistent fixed point 0.

This kind of bifurcation is called “period-doubling”. An example is that of the family

$$f_\lambda(x) = -x + x^2 + \lambda x$$

Another example occurs in the quadratic family $f_\lambda(x) = \lambda x(1 - x)$, when we pass the value $\lambda = 3$ of the parameter. You may check this with my applet <http://w3.math.uminho.pt/~scosentino/bestiario/logistic.html>.

Doing simulations with a computer, Mitchell J. Feigenbaum discovered, in the '70 of the last century, that certain families of transformations produce a “cascade” of period-doublings, i.e. there is a sequence of values $\lambda_1 < \lambda_2 < \dots < \lambda_n < \lambda_{n+1} \dots$ of the parameter such that, when passing through λ_{n+1} orbits of period 2^{n+1} are created in a neighbourhood of orbits of period 2^n , created by the previous value λ_n . This phenomenon is easily observed, and indeed seems to be “universal”: it happens for almost all families, provided we find the region where it takes place. The following picture is obtained if we plot the parameter λ , within some interval, versus a typical orbit of f_λ , say $\{x_{100}, x_{101}, \dots, x_{200}\}$ starting from a random initial point x_0 . Here is what you get.



Bifurcation diagram for the logistic family.

(from https://en.wikipedia.org/wiki/Period-doubling_bifurcation)

Even more mysterious is that, as already observed by Feigenbaum, the limit $\lambda_\infty = \lim_{n \rightarrow \infty} \lambda_n$ seems to exist, it is achieved exponentially, i.e. $|\lambda_\infty - \lambda_n| \simeq \text{const} \times \delta^{-n}$ where

$$\delta = \lim_{n \rightarrow \infty} \frac{\lambda_n - \lambda_{n-1}}{\lambda_{n+1} - \lambda_n} \simeq 4.669201609102990671853 \dots$$

seem to be independent from the family! This mystery was explained later by Lanford, Epstein, Sullivan, ...

7 Statistical description of orbits

Together with the topological point of view, a source of informations about dynamical systems is their statistical description. The idea is to measure the relative size of those points whose orbits have certain definite properties. This is done looking for invariant probability measures, and the main result is the Birkhoff-Khinchin ergodic theorem. To state and prove the Birkhoff-Khinchin ergodic theorem, we need to recall many standard facts and results of integration theory. You can find most of them in the classical [Ru87] or [Ha74].

7.1 Probability measures

Probability spaces. A *measurable space* is a pair (X, \mathcal{E}) , a non-empty set X together with a σ -algebra of subsets \mathcal{E} . Recall that a (Boolean) algebra is a nonempty family \mathcal{A} of subsets of X which contains X , which contains the complement of any of its elements, and which is closed under finite unions and intersections. A σ -algebra is an algebra which is also closed under countable unions and intersections. Given any family \mathcal{C} of subsets of X , there exists a minimal σ -algebra $\sigma(\mathcal{C})$ which contains all the elements of \mathcal{C} , which is called the σ -algebra generated by \mathcal{C} .

If (X, τ) is a topological space, the *Borel σ -algebra* is $\sigma(\tau)$, the smallest σ -algebra which contains all open sets.

A *measure* on the measurable space (X, \mathcal{E}) is a σ -additive function $\mu : \mathcal{E} \rightarrow [0, \infty]$ such that $\mu(\emptyset) = 0$. Here σ -additivity means that, if (S_n) is a countable family of pairwise disjoint elements of \mathcal{E} , then

$$\mu(\cup_n S_n) = \sum_n \mu(S_n)$$

The triple $(\Omega, \mathcal{E}, \mu)$ is said a *measure space*, or *probability space* if it happens that $\mu(X) = 1$. Given a probability space, measurable sets $A \in \mathcal{E}$ are commonly called "events", and the number $\mu(A)$ is called "probability of the event A ". Basic properties of probability measures are the following: probability measures are monotone, i.e. $\mu(S) \leq \mu(T)$ if $S \subset T$, and σ -subadditive, i.e. if (S_n) is a countable family of elements of \mathcal{E} then

$$\mu(\cup_n S_n) \leq \sum_n \mu(S_n)$$

Probability measures are continuous from below and from above, in the following sense: if $S_n \uparrow S$ then $\mu(S_n) \uparrow \mu(S)$, and if $S_n \downarrow S$ then $\mu(S_n) \downarrow \mu(S)$. Both continuity properties are equivalent, and indeed a simple argument shows that they are equivalent to continuity from above at \emptyset : if $S_n \downarrow \emptyset$ then $\mu(S_n) \downarrow 0$. Moreover, continuity is equivalent to σ -additivity if the set function μ is only assumed (finitely) additive.

A subset $E \subset X$ has *zero measure* if it is contained in a measurable set $S \in \mathcal{E}$ with $\mu(S) = 0$. If any set with zero measure belongs to \mathcal{E} , then the measure space (X, \mathcal{E}, μ) is said *complete*. Any measure space can be canonically completed, extending the measure to the σ -algebra $\bar{\mathcal{E}}$ made of \mathcal{E} and of subsets of zero measure. A property (like continuity of a function, or convergence of a sequence of functions) holds μ -a.e. ("almost everywhere" with respect to the measure μ) if the set of points of X where it does not hold has zero measure.

Construction of probability measures. Measures are never "explicitly" given as functions on a σ -algebra. A set function $\mu : \mathcal{P}(X) \rightarrow [0, \infty]$ is an *exterior measure* if it is monotone, σ -subadditive, and if $\mu(\emptyset) = 0$. It happens that, given an exterior measure μ , the family of μ -measurable sets, defined as

$$\mathcal{E} = \{E \subset X \text{ such that } \mu(S) = \mu(S \cap E) + \mu(S \cap E^c) \text{ for any } S \subset X\}$$

is a σ -algebra, and that μ is a complete measure if restricted on \mathcal{E} (the proof is quite long and delicate, but the only idea it uses is the following: in order to check that $E \in \mathcal{E}$ it is indeed sufficient, by virtue of monotonicity and subadditivity of μ , to check that $\mu(S) \geq \mu(S \cap E) + \mu(S \cap E^c)$ for any $S \subset X$). A strategy to construct interesting measures on uncountable spaces is: start with an exterior measure (it is very easy to produce exterior measures, for example by means of variational principles) and then check that the σ -algebra of measurable sets is sufficiently big for our purpose.

The idea of Carathéodory is the following. A *probability measure* on an algebra \mathcal{A} of subsets of X is an additive function $m : \mathcal{A} \rightarrow [0, 1]$ such that $m(\emptyset) = 0$, $m(X) = 1$, and such that $A_n \downarrow \emptyset$ implies $m(A_n) \downarrow 0$. Given a probability measure m on an algebra \mathcal{A} , the recipe

$$\mu(S) = \inf \left\{ \sum m(A_n) \text{ with } S \subset \cup_n A_n \text{ e } A_n \in \mathcal{A} \right\}$$

defines an exterior measure on $\mathcal{P}(X)$, hence the above construction produces a measure μ on the σ -algebra of μ -measurable sets, which contains \mathcal{A} and so contains $\sigma(\mathcal{A})$. One then checks that $\mu(A) = m(A)$ for any $A \in \mathcal{A}$, so that μ is an “extension” of the measure m . Carathéodory’s extension theorem is then stated in the following form:

Theorem 7.1 (Carathéodory’s extension theorem). *Given a probability measure m on a algebra \mathcal{A} of subsets of X , there exists a unique measure μ on $\sigma(\mathcal{A})$ which extends m .*

The following corollary of Carathéodory’s theorem is also useful, for example when trying to prove that some event has a definite probability.

Theorem 7.2 (Approximation theorem). *Let (X, \mathcal{E}, μ) be a probability space, and let \mathcal{A} be an algebra of subsets of X such that $\sigma(\mathcal{A}) = \mathcal{E}$. Then, for any $A \in \mathcal{E}$ and any $\varepsilon > 0$, we can find a $A_\varepsilon \in \mathcal{A}$ such that*

$$\mu(A \Delta A_\varepsilon) < \varepsilon.$$

Indeed, one easily sees that the family $\mathcal{C} = \{A \in \mathcal{E} \text{ s.t. } \forall \varepsilon > 0 \exists A_\varepsilon \in \mathcal{A} \text{ s.t. } \mu(A \Delta A_\varepsilon) < \varepsilon\}$ is a σ -algebra. Since \mathcal{A} is obviously contained in \mathcal{C} , this implies that $\mathcal{E} = \sigma(\mathcal{A}) \subset \mathcal{C} \subset \mathcal{E}$.

Lebesgue measure. The collection \mathcal{I} of intervals of the real line is a *semi-algebra*, i.e. the intersection of two elements of \mathcal{I} is in \mathcal{I} and the complement of an element of \mathcal{I} is a union of elements of \mathcal{I} . The function $m : \mathcal{I} \rightarrow [0, \infty]$, defined as $m([a, b]) = |b - a|$ if a e b are finite, and ∞ if the interval is unbounded, is monotone and gives value zero to the empty set. Postulating additivity, the function m extends to a measure on the algebra \mathcal{A} made of disjoint unions of elements of \mathcal{I} (this is not trivial!, the proof uses the Heine-Borel theorem about compact subsets of the real line). The function $\mu : \mathcal{P}(\mathbb{R}) \rightarrow [0, \infty]$, defined as

$$\mu(E) = \inf \left\{ \sum m(C_n) \text{ with } E \subset \cup_n C_n \text{ e } C_n \in \mathcal{A} \right\}$$

is then an exterior measure on the real line. The σ -algebra \mathcal{L} of μ -measurable sets, called *Lebesgue σ -algebra*, contains the Borel sets, because it contains the intervals. The restriction $\ell = \mu|_{\mathcal{L}}$, as well as $\mu|_{\mathcal{B}(\mathbb{R})}$, is called *Lebesgue measure*.

Observe that Lebesgue measure on the real line is not a probability measure, having infinite mass. Nevertheless, one can easily define probability measures on bounded intervals taking normalized restrictions of Lebesgue measure. For example, take $X = [0, 1]$, and $\mathcal{E} = \mathcal{B}(X) = \{X \cap B \text{ with } B \in \mathcal{B}(\mathbb{R})\}$, the Borel subsets of the interval. The restriction of ℓ to \mathcal{E} is a probability measure, called Lebesgue measure on the unit interval.

The very same construction works in \mathbb{R}^N , starting with the semi-algebra of “rectangles” measured by the “euclidean volume”, and produces a measure ℓ on $\mathcal{B}(\mathbb{R}^N)$, also called Lebesgue measure. Lebesgue measure is the unique measure over the Borel sets of the euclidean space which is invariant under translations, i.e. $\ell(\lambda + B) = \ell(B)$ for any $\lambda \in \mathbb{R}^N$ and any Borel set B , and which is normalized to give measure one to the unit square, i.e. $\ell([0, 1]^N) = 1$.

The axiom of choice allows one to “give examples” of subsets which are not Lebesgue-measurable (for example, the set made of one point for each orbit of an irrational rotation of the circle).

The following result is useful (see [Mat95] for a proof). Below, $B_\varepsilon(x) = \{y \in \mathbb{R}^N \text{ s.t. } \|x - y\|_2 < \varepsilon\}$ denotes the open ball of radius $\varepsilon > 0$ and center $x \in \mathbb{R}^N$ w.r.t. the Euclidean distance $\|x - y\|_2^2 = \sum_{i=1}^N |x_i - y_i|^2$.

Theorem 7.3 (Lebesgue density theorem). *Let $A \subset \mathbb{R}^N$ be a Lebesgue-measurable set. For ℓ -almost any $x \in A$ there exists the limit*

$$\lim_{\varepsilon \rightarrow 0} \frac{\ell(A \cap B_\varepsilon(x))}{\ell(B_\varepsilon(x))} = 1$$

Kolmogorov extension. Let X be a finite space, equipped with the discrete topology, and let Σ^+ be the topological product $X^{\mathbb{N}} = \{x : \mathbb{N} \rightarrow X\}$, its point identified with sequences $x = (x_1, x_2, \dots, x_n, \dots)$ with $x_n \in X$. Let \mathcal{C} be the collection of *cylinders* of X , the subsets of the form

$$C_B = \{x \in \Sigma^+ \text{ s.t. } (x_1, x_2, \dots, x_n) \in B\}$$

with B an open subset of X^n . Cylinders form a basis of the product topology of Σ^+ , which makes Σ^+ a compact metrizable space. In particular, the Borel σ -algebra of Σ^+ is $\mathcal{B} = \sigma(\mathcal{C})$. Let $\mu_1, \mu_2, \mu_3, \dots, \mu_n, \dots$ be probability measures defined on the Borel sets of $X, X^2, \dots, X^n, \dots$, respectively. The sequence (μ_n) is said *consistent* if

$$\mu_{n+1}(B \times X) = \mu_n(B)$$

for any n and any Borel subset $B \subset X^n$. The (most elementary version of) Kolmogorov extension theorem says that

Theorem 7.4 (Kolmogorov extension theorem). *Given a consistent family of probability measures as above, there exists a unique probability measure μ , defined on the Borel σ -algebra of Σ^+ , such that*

$$\mu(C_B) = \mu_n(B)$$

for any cylinder C_B .

Proof. The proof consists in the following two steps. First, observe that cylinders form an algebra, and use consistency of the μ_n 's to verify that the formula above does define a function $\mu : \mathcal{C} \rightarrow [0, 1]$ on cylinders (i.e. it does not depend on the different ways the same cylinder may be presented) which is additive and properly normalized. Then, use compactness of X to check that μ is continuous at \emptyset , in order to apply Carathéodory theorem. Indeed, let (A_n) be a sequence of cylinders such that $A_n \downarrow \emptyset$, and assume by contradiction that $\mu(A_n) > \delta > 0$ for any n . This implies that $A_n \neq \emptyset$ for any n , but, since the A_n are compact, then the Cantor intersection theorem says that $\bigcap_n A_n \neq \emptyset$, contrary to the hypothesis. \square

Kolmogorov theorem is the key tool in probability theory, since it allows one to construct measures which describe an infinite sequence of trials starting with some rule which gives information about the n -th trial given the knowledge of the first $n-1$. It actually works with much more general spaces and in a more general setting. Also, one can easily adapt the construction to $\prod_{n \in \mathbb{N}} X_n$, the topological product of a countable family of finite spaces. In some precise sense, this is a universal model of a dynamical system.

Bernoulli trials. If $X = \{0, 1\}$, then $\Sigma^+ = X^{\mathbb{N}}$ is the state space of infinite Bernoulli trials with two possible outcomes: success and failure. Let $\mu_1 : \mathcal{P}(X) \rightarrow [0, 1]$ be a any probability measure, defined by $\mu_1(\{1\}) = p$. Kolmogorov construction can be applied postulating the independence of different trials, i.e. declaring that the family formed by the cylinders $\{x_n = 1\}$ is an independent family, and giving measure p to each $\{x_n = 1\}$. The resulting probability space $(\Sigma^+, \mathcal{B}, \mu)$ describes the infinite independent Bernoulli trials. Of course, the very same construction can be made when X is a finite space with any finite number z of elements.

7.2 Transformations and invariant measures

Measurable transformations. A transformation $f : X \rightarrow X$ of the measurable space (X, \mathcal{E}) is said *measurable* if $f^{-1}(A) \in \mathcal{E}$ for any $A \in \mathcal{E}$. A measurable transformation f is said an *endomorphism* of the measurable space, or an *automorphism* if it is invertible and its inverse is measurable too.

Observe that an endomorphism f of a measurable space (X, \mathcal{E}) acts naturally on the space of measures on \mathcal{E} by "push forward": if μ is a measure, then $f_*\mu$, defined by $(f_*\mu)(A) = \mu(f^{-1}(A))$ for any $A \in \mathcal{E}$, is also a measure.

Let f be an endomorphism of the measurable space (X, \mathcal{E}) . A probability measure μ on \mathcal{E} is *invariant* (w.r.t. the transformation f) if $f_*\mu = \mu$, namely if

$$\mu(f^{-1}(A)) = \mu(A)$$

for any $A \in \mathcal{E}$. If this happens, we also say that f is an *endomorphism* (resp. an *automorphism*) of the probability space (X, \mathcal{E}, μ) . The meaning of this definition is that "mean values" of integrable observables $\varphi : X \rightarrow \mathbb{R}$ with respect to invariant probability measures do not change with time, in the sense that $\int_X \varphi d\mu = \int_X (\varphi \circ f) d\mu$.

Given an endomorphism f of the probability space (X, \mathcal{E}, μ) , one says that an event $A \in \mathcal{E}$ is *invariant mod 0* if $\mu(A \Delta f^{-1}(A)) = 0$. The set of invariant mod 0 events form a sub- σ -algebra of \mathcal{E} , denoted by \mathcal{E}_f .

How to prove that a measure is invariant. The very definition of invariance does not help too much if we want to prove that a certain measure μ on the σ -algebra \mathcal{E} is invariant w.r.t. the measurable transformation $f : X \rightarrow X$. The trick is the following. Suppose that we can prove that $\mu(f^{-1}(C)) = \mu(C)$ for any $C \in \mathcal{C}$, where \mathcal{C} is some subset of \mathcal{E} . Caratheodory theorem implies that $f_*\mu$ and μ are the same measure on the σ -algebra $\sigma(\mathcal{C})$ generated by \mathcal{C} . On the other side, the family of those $A \in \mathcal{E}$ such that $\mu(f^{-1}(A)) = \mu(A)$ is easily seen to be a σ -algebra. Hence, if it happens that $\sigma(\mathcal{C}) = \mathcal{E}$, then μ is actually invariant. In other words, in order to prove that μ is invariant it is sufficient to check that $\mu(f^{-1}(C)) = \mu(C)$ for any C belonging to a family of subsets of X which generate the σ -algebra \mathcal{E} .

Observables as random variables. When dealing with a endomorphism $f : X \rightarrow X$ of the probability space (X, \mathcal{E}, μ) , one should consider *measurable* observables $\varphi : X \rightarrow \mathbb{R}$ (or \mathbb{C}), those functions such that $\varphi^{-1}(A) \in \mathcal{E}$ for any Borel set $A \subset \mathbb{R}$. In the context of probability theory they are called "random variables", and the sequence of observables $\varphi \circ f^n$ may be interpreted as a "random process". If φ is integrable, the Lebesgue integral $\int_X \varphi d\mu$ is interpreted as the "mean value" of φ . Of course, invariance of a measurable observable must be intended modulo sets of zero measure. Then, one can consider the Banach spaces $L^p(\mu)$ of (equivalence classes of) observables equipped with the L^p -norm

$$\|\varphi\|_p = \left(\int_X |\varphi|^p d\mu \right)^{1/p}$$

and use the full power of integration theory to get informations about the dynamical system. In particular, $L^2(\mu)$ is a Hilbert space if equipped with the inner product

$$\langle \varphi, \psi \rangle = \int_X \varphi \bar{\psi} d\mu$$

Conditional mean. Recall that, given a measurable space (X, \mathcal{E}) , a measure ν is said *absolutely continuous* w.r.t. the measure μ if $\nu(A) = 0$ whenever $\mu(A) = 0$. The following technical result (which may be proved using Hilbert space techniques) is particularly useful:

Theorem 7.5 (Radon-Nikodym). *Let (X, \mathcal{E}, μ) be a probability space, and let ν be a finite measure over \mathcal{E} which is absolutely continuous with respect to μ . Then there exists a nonnegative integrable random variable ρ (called the Radon Nikodym derivative of ν w.r.t. μ and denoted by $d\nu/d\mu$) such that*

$$\nu(A) = \int_A \rho d\mu$$

for any $A \in \mathcal{E}$.

A particularly important tool, taken from the theory of probability, is the conditional mean. Let (X, \mathcal{E}, μ) be a probability space, and let \mathcal{F} be a sub- σ -algebra of \mathcal{E} . Given an integrable random variable φ , there exists a unique random variable $\varphi_{\mathcal{F}}$, called the *conditional mean* of φ w.r.t. \mathcal{F} , which is \mathcal{F} -measurable (i.e. the inverse image of any Borel set belongs to \mathcal{F}) and such that

$$\int_A \varphi_{\mathcal{F}} d\mu = \int_A \varphi d\mu$$

for any $A \in \mathcal{F}$. Indeed, if $\varphi \geq 0$, then one can define $\varphi_{\mathcal{F}}$ as equal to the Radon-Nikodym derivative of the measure $A \mapsto \int_A \varphi d\mu$, defined on \mathcal{F} , with respect to the restriction $\mu|_{\mathcal{F}}$. The general case is treated by linearity, writing φ as a difference of two non-negative random variables. Uniqueness is intended μ -a.e., i.e. modulo sets of zero probability. The conditional mean is monotone, namely if $\varphi \geq 0$ then $\varphi_{\mathcal{F}} \geq 0$, and preserves the mean value, since $\int_X \varphi_{\mathcal{F}} d\mu = \int_X \varphi d\mu$. It can be considered as a "projection" of φ onto the space of \mathcal{F} -measurable random variable, preserving the mean value. In particular, if \mathcal{N} is the trivial σ -algebra made of events of measure 0 or 1, then $\varphi_{\mathcal{N}}$ is constant a.e. and equal to $\int_X \varphi d\mu$.

Topological dynamical systems and Borel measures. If we are interested in the dynamics of a continuous transformation $f : X \rightarrow X$ of a topological space X , it is natural to consider the Borel σ -algebra \mathcal{B} , the smallest σ -algebra of subsets of X which contain all open sets. The map f is then an endomorphism of (X, \mathcal{B}) . Probability measures on \mathcal{B} are said *Borel probability measures*. If, moreover, X is a compact metric space, one can consider the space $\mathcal{C}^0(X, \mathbb{R})$ of bounded continuous real valued functions of X (observe that, since X is compact, any continuous function is automatically bounded), equipped with the sup norm

$$\|\varphi - \psi\|_{\infty} = \sup_{x \in X} |\varphi(x) - \psi(x)|$$

These observables are clearly integrable w.r.t. to any Borel probability measure μ , and the mean value map

$$\varphi \mapsto \int_X \varphi d\mu$$

is a bounded, positive definite (in the sense that $\int_X \varphi d\mu \geq 0$ for any $\varphi \geq 0$) linear functional on $\mathcal{C}^0(X, \mathbb{R})$. The basic fact about Borel measures is the converse of that, namely

Theorem 7.6 (Riesz-Markov representation theorem). *Let X be a compact metric space. Given any bounded and positive definite linear functional L on $\mathcal{C}^0(X, \mathbb{R})$ such that $L(1) = 1$, there exists a unique Borel probability measures μ such that*

$$L(\varphi) = \int_X \varphi d\mu$$

for any $\varphi \in \mathcal{C}^0(X, \mathbb{R})$

The space of invariant probability measures. The space Prob of probability measures on a measurable space (X, \mathcal{E}) has a natural convex structure: convex combinations of probability measures are also probability measures. An arbitrary measurable transformation $f : X \rightarrow X$ of a measurable space may not admit any invariant probability measure. On the other side, if μ_0 and μ_1 are invariant probability measures, so are their convex combinations $\mu_t = (1-t)\mu_0 + t\mu_1$ for any $t \in [0, 1]$. This means that the set Prob_f of invariant probability measures on \mathcal{E} is a convex set: if it contains two points, it contains the whole segment between them.

Now, let (X, d) be a compact metric space and let \mathcal{B} its Borel σ -algebra. The space Prob of probability measures on \mathcal{B} can be equipped with a natural topology, called the *weak* topology*, which says essentially that two measures are near if they give nearby mean values to some well behaved observables. Formally, one says that a sequence of measures (μ_n) converge weakly* to a measure μ , which we denote simply as $\mu_n \rightarrow \mu$, if

$$\int_X \varphi d\mu_n \rightarrow \int_X \varphi d\mu$$

for any (bounded) continuous function $\varphi : X \rightarrow \mathbb{R}$. The space $\mathcal{C}^0(X, \mathbb{R})$ of bounded continuous real valued functions on X , equipped with the sup norm, is a separable Banach space. In particular, it admits a countable set of points $\{\varphi_n\}_{n \in \mathbb{N}}$ which is dense in its unit sphere. Given that, one defines, for any couple of Borel probability measures μ and ν , a distance

$$d(\mu, \nu) = \sum_{n=1}^{\infty} 2^{-n} \cdot \left| \int_X \varphi_n d\mu - \int_X \varphi_n d\nu \right|$$

It turns out that d is indeed a metric, and that it induces the weak* topology on Prob. The important fact (somewhere called "Helly's theorem"), which follows from the Ascoli-Arzelà theorem together with the above Riesz-Markov representation theorem, is that Prob, equipped with the weak* topology, is a compact space: any sequence (μ_n) of Borel probability measures admits a weakly* convergent subsequence $\mu_{n_i} \rightarrow \mu$.

Now, we are in position to prove the existence of invariant probability measures for certain well behaved dynamical systems.

Theorem 7.7 (Krylov-Bogolyubov). *A continuous transformation $f : X \rightarrow X$ of a metrizable compact space X admits at least one Borel invariant probability measure.*

Proof. Take any Borel probability measure μ_0 on X , and inductively define a family of probability measures μ_n by $\mu_{n+1} = f_*\mu_n$. Consider the family of Cesaro means

$$\bar{\mu}_n = \frac{1}{n+1} \sum_{k=0}^n \mu_k$$

Since the space of Borel probability measures on a compact metrizable space is compact w.r.t. weak* convergence, there exist a weakly* convergent subsequence $\bar{\mu}_{n_i} \rightarrow \mu$. One then easily sees that

$$\begin{aligned} \int_X (\varphi \circ f) d\mu &= \lim_{i \rightarrow \infty} \frac{1}{n_i + 1} \sum_{k=0}^{n_i} \int_X (\varphi \circ f) d\mu_k \\ &= \lim_{i \rightarrow \infty} \frac{1}{n_i + 1} \sum_{k=0}^{n_i} \int_X \varphi d\mu_{k+1} \\ &= \lim_{i \rightarrow \infty} \frac{1}{n_i + 1} \sum_{k=0}^{n_i} \int_X \varphi d\mu_k + \frac{1}{n_i + 1} \left(\int_X \varphi d\mu_{n_i+1} - \int_X \varphi d\mu_0 \right) \\ &= \int_X \varphi d\mu \end{aligned}$$

for any bounded continuous observable φ , hence that μ is an invariant measure. \square

7.3 Invariant measures and time averages

The relevance of invariant measures when studying the dynamics of continuous transformations is due to the following crucial observations.

Invariant measures and time averages. Assume that, for a given point $x \in X$, the time averages

$$\bar{\varphi}(x) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n \varphi(f^k(x))$$

do exist for any bounded continuous observable φ . One easily shows that the functional $\mathcal{C}_b^0(X, \mathbb{R}) \rightarrow \mathbb{R}$ defined by $\varphi \mapsto \bar{\varphi}(x)$ is linear, bounded and positive definite. There follows from the Riesz-Markov representation theorem that there exists a unique Borel probability measure μ_x on X such that

$$\bar{\varphi}(x) = \int_X \varphi d\mu_x$$

for any $\varphi \in \mathcal{C}_b^0(X, \mathbb{R})$. The invariance property $\bar{\varphi}(x) = (\bar{\varphi} \circ f)(x)$ for time averages then implies that $\int_X (\varphi \circ f) d\mu_x = \int_X \varphi d\mu_x$ for any φ , hence that μ_x is an invariant probability measure. In the language of physicists, this says that "time averages" along the orbit of x are equal to "space averages" with respect to the measure μ_x .

One is thus lead to consider the following questions. Do there exist points x for which time averages exists? Given an invariant measure μ , do there exist, and how many, points x such that $\mu = \mu_x$?

Periodic orbits. Here is a trivial but important example. Let p be a periodic point with period n . The time average $\bar{\varphi}(p)$ of any observable φ exists, and is equal to the arithmetic mean of φ along the orbit, namely

$$\bar{\varphi}(p) = \frac{1}{n} \sum_{k=0}^{n-1} \varphi(f^k(p))$$

If μ_p denotes the normalized sum $\frac{1}{n} \sum_{k=0}^{n-1} \delta_{f^k(p)}$ of Dirac masses placed on the orbit of p , this amounts to say that $\bar{\varphi}(p) = \int_X \varphi d\mu_p$.

Let p be a fixed point, and $\varphi : X \rightarrow \mathbb{R}$ be an observable which is continuous at p . If $x \in W^s(p)$, then the time average $\bar{\varphi}(x)$ exists and is equal to $\varphi(p)$, i.e. time averages of points in the basin of attraction of p are described by the Dirac measure $\mu_p = \delta_p$.

The Birkhoff-Khinchin ergodic theorem. Ergodic theorems are the milestones of ergodic theory, and deal with various type of convergence of the time means $\bar{\varphi}_n$ for certain classes of observables φ . In particular, the Birkhoff-Khinchin ergodic theorem³¹ must be thought as the generalization of the Kolmogorov strong law of large numbers, as it says that time means of certain well-behaved observables exist almost everywhere. The Birkhoff-Khinchin ergodic theorem was actually preceeded by the von Neumann's "statistic" ergodic theorem³²

If $f : X \rightarrow X$ is an endomorphism of the probability space (X, \mathcal{E}, μ) , one can consider the "shift" operator $U_f : L^2(\mu) \rightarrow L^2(\mu)$, defined by $(U_f \varphi)(x) := \varphi(f(x))$. It is clearly an isometry, and if f is a homeomorphism it is unitary. The fixed point set of U_f is the space of invariant L^2 -observable. The von Neumann theorem then asserts convergence of time means $\bar{\varphi}_n \rightarrow \bar{\varphi}$ in the Hilbert space $L^2(\mu)$.

Theorem 7.8 (von Neumann "statistic" ergodic theorem). *Let U be a unitary operator on a Hilbert space \mathcal{H} , let $\mathcal{H}_U = \{v \in \mathcal{H} \text{ s.t. } Uv = v\}$ denote the closed subspace of those vectors which are fixed by U , and $P_U : \mathcal{H} \rightarrow \mathcal{H}_U$ denote the orthogonal projection onto \mathcal{H}_U . Then, for any vector $\varphi \in \mathcal{H}$ we have*

$$\lim_{n \rightarrow \infty} \left\| \frac{1}{n+1} \sum_{k=0}^n U^k \varphi - P_U \varphi \right\|_{\mathcal{H}} = 0$$

Here, we prove the stronger

Theorem 7.9 (Birkhoff-Khinchin "individual" ergodic theorem). *Let $f : X \rightarrow X$ be an endomorphism of the probability space (X, \mathcal{E}, μ) , and let $\varphi \in L^1(\mu)$ be an integrable observable. Then the limit*

$$\bar{\varphi}(x) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n \varphi(f^k(x))$$

exists for μ -almost any $x \in X$. Moreover, the observable $\bar{\varphi}$ is in $L^1(\mu)$, is invariant, and satisfies

$$\int \bar{\varphi} d\mu = \int \varphi d\mu$$

Proof. (by A. Garsia, as explained in [KH95]) Let \mathcal{E}_f be the invariant σ -algebra. For any $\psi \in L^1$, set $\psi_n = \max_{k \leq n} \sum_{k=0}^{n-1} \varphi \circ f^k$ and observe that $E_\psi = \{x \in X \text{ s.t. } \psi_n(x) \rightarrow \infty\} \in \mathcal{E}_f$. One easily sees that the sequence $\psi_{n+1} - \psi_n \circ f$ is decreasing, and converges to ψ at the points of E_ψ . The monotone convergence theorem and the invariance of μ imply that

$$0 \leq \int_{E_\psi} (\psi_{n+1} - \psi_n) d\mu = \int_{E_\psi} (\psi_{n+1} - \psi_n \circ f) d\mu \rightarrow \int_{E_\psi} \psi d\mu = \int_{E_\psi} \psi_{\mathcal{E}_f} d\mu|_{\mathcal{E}_f}$$

In particular, if $\psi_{\mathcal{E}_f} < -\varepsilon < 0$ then $\mu(E_\psi) = 0$. On the other side,

$$\limsup \frac{1}{n} \sum_{k=0}^{n-1} \psi \circ f^k(x) \leq \limsup \frac{1}{n} \psi_n \leq 0$$

³¹G.D. Birkhoff, Proof of the ergodic theorem, *Proc. Natl. Acad. Sci. USA* **17** (1931) 656-660.

³²J. von Neumann, Proof of the Quasi-ergodic Hypothesis, *Proc. Natl. Acad. Sci. USA* **18** (1932), 70-82.

on $X \setminus E_\psi$. Applying twice these observations to the observables $\varphi - \varphi_{\mathcal{E}_f} - \varepsilon$ and $-\varphi + \varphi_{\mathcal{E}_f} - \varepsilon$, with $\varepsilon > 0$, we find

$$\limsup \frac{1}{n} \sum_{k=0}^{n-1} \varphi \circ f^k(x) - \varphi_{\mathcal{E}_f} - \varepsilon \leq 0 \quad \liminf \frac{1}{n} \sum_{k=0}^{n-1} \varphi \circ f^k(x) - \varphi_{\mathcal{E}_f} + \varepsilon \geq 0$$

μ -almost everywhere. Since ε was arbitrary, the limit $\overline{\varphi}(x)$ exists and is equal to $\varphi_{\mathcal{E}_f}(x)$ for μ -almost every x . The rest of the theorem then follows easily from the properties of the conditional mean. \square

7.4 Examples of invariant measures

Haar measures. Any locally compact topological group G admits a *Haar measure*, a measure μ on its Borel sets which is left-invariant, i.e. satisfies $L_g \mu = \mu$ for any $g \in G$. Moreover, the Haar measure is unique up to a constant factor. It is an exercise that μ is a finite measure, hence can be renormalized to give a probability measure, iff G is compact. There follows that translations on compact topological groups admits invariant probability measures.

On the other side, for some groups G , called *unimodular*, the Haar measure μ is both left and right invariant. If $\Gamma \subset G$ is a lattice, i.e. a subgroup such that $\mu(G/\Gamma) < \infty$, then the normalized Haar measure on the homogeneous space G/Γ is an invariant probability measure for any left translation $g\Gamma \mapsto sg\Gamma$.

Rotations of the circle. Lebesgue probability measure ℓ on the circle is invariant for the rotations $R_\alpha : x + \mathbb{Z} \mapsto x + \alpha + \mathbb{Z}$, with $\alpha \in \mathbb{R}$. Indeed, rotations of the circle are isometries, and the Lebesgue measure $\ell(I)$ of an interval is its "length".

Coverings of the circle. Lebesgue probability measure ℓ on the circle is invariant for the maps $E_N : x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$, with $N \in \mathbb{Z} \setminus \{0\}$. This comes from the fact that the inverse image of a sufficiently small interval I with length $\ell(I)$ is the disjoint union of $|N|$ intervals with length $\ell(I)/|N|$.

Bernoulli shifts. Consider the Bernoulli shift $\sigma : \Sigma^+ \rightarrow \Sigma^+$ over the alphabet $X = \{1, 2, \dots, N\}$. Let p be a "probability on X ", i.e. a finite set of nonnegative numbers p_1, p_2, \dots, p_N such that $p_1 + p_2 + \dots + p_N = 1$. Given a centered cylinder C_α , we define $\mu(C_\alpha)$ as equal to the product $p_{\alpha_1} p_{\alpha_2} \dots p_{\alpha_n}$. This function μ extends in a unique way as a finitely additive function on the algebra \mathcal{A} generated by the centered cylinders, the algebra which contains all finite unions of centered cylinders as well as the empty set and Σ^+ . One then show that μ is σ -additive on \mathcal{A} (for example, showing that if a decreasing sequence $A_1 \supset A_2 \supset \dots$ has empty intersection then $\mu(A_n) \rightarrow 0$). Since centered cylinders generates the topology of Σ^+ , Carathéodory theorem implies that there exists a unique extension, which we still call μ , of this measure on the Borel σ -algebra of Σ^+ . This measure is called the *Bernoulli measure* defined by p .

As for the "physical" meaning of this measure, you may imagine that X represents the possible outcomes when tossing a coin with z sides, and p_k is the probability of obtaining the k -th side. Then points in Σ^+ represent the outcomes of an infinite sequence of tossings, and the very definition of μ says that each trial is described by the probability p , and each trial is "independent" from any finite collection of different trials.

It is not surprising that μ is indeed an invariant probability measure. This comes from the fact that the inverse image $\sigma^{-1}(A)$ of any $A \in \mathcal{A}$ is the disjoint union of N elements B_1, B_2, \dots, B_N of the algebra (obtained from A choosing the first letter in z different ways) with measures $\mu(B_k) = p_k \cdot \mu(A)$, so that

$$\mu(\sigma^{-1}(A)) = \sum_{k=1}^N p_k \cdot \mu(A) = \mu(A)$$

Absolutely continuous invariant measures for maps and flows. Let U be a domain in some euclidean \mathbb{R}^N , and let ℓ denote the Lebesgue measure on U , given locally as

$$d\ell = dx = dx_1 dx_2 \dots dx_n.$$

Thus, the volume of an open set $A \subset U$ is $\ell(A) = \int_A dx$. A local diffeomorphism $f : U \rightarrow U$ of class C^1 preserves the measure vol iff

$$\sum_{x \in f^{-1}\{x'\}} \frac{1}{|\det f'(x)|} = 1$$

for any point $x' \in U$, as one can check using the change of coordinates formula. Also interesting is to see wheather f preserves an absolutely continuous measure $d\mu = \rho d\ell$, defined by $\mu(A) = \int_A \rho dx$. This happens iff the "density" ρ satisfies the equation

$$\sum_{x \in f^{-1}\{x'\}} \frac{\rho(x)}{|\det f'(x)|} = \rho(x')$$

for any point $x' \in U$.

Now, let Φ_t be the flow of a vector field $v = \sum_{k=1}^n v_k \frac{\partial}{\partial x_k}$ on U . The above obviously applies, considering the Jacobian of the diffeomorphisms Φ_t . Since

$$\det \Phi'_t = \int_0^t \text{div} v \circ \Phi_s ds$$

we get the result that Lebesgue measure ℓ is invariant under the flow of v iff

$$\text{div} v = \sum_{k=1}^n \frac{\partial v_k}{\partial x_k} = 0$$

In general, the absolutely continuous measure $d\mu = \rho d\ell$ is invariant under the flow of v iff its density satisfies $\text{div}(\rho v) = 0$.

Hamiltonian flows. Consider a symplectic manifold (X, ω) . Liouville measure $d\mu = \omega^n$ is invariant under the Hamiltonian flow of any Hamiltonian function H . If X has finite volume, it can be normalized to give an invariant probability measure.

Geodesic flows. Consider a geodesic flow on the unit tangent bundle $\pi : SM \rightarrow M$ of the Riemannian manifold (M, g) . Let $d\mu = \sqrt{g} dx$ denote the Riemannian volume form on M , and let $d\sigma_m$ denotes the Lebesgue probability measure on the sphere $S_m M = \pi^{-1}\{m\}$. The Liouville measure ℓ , defined locally as $d\mu(m) \times d\sigma_m$, is invariant under the geodesic flow.

Gauss map. Any irrational real number $x \in (0, 1]$ has a unique continued fraction representation of the form

$$x = [0; a_1, a_2, a_3, \dots] = \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3 + \frac{1}{\ddots}}}}$$

where the a_n are nonnegative integers. The equality sign and the "infinite fraction" above mean that the sequence of finite continued fractions

$$p_n/q_n = [0; a_1, a_2, \dots, a_n] = \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\ddots + \frac{1}{a_n}}}}$$

which are called “convergents”, do converge to x as $n \rightarrow \infty$. The sequence of partial quotients a_n is inductively constructed as follows. First, observe that if $a_1 = [1/x]$ and $x_1 = 1/x - a_1$ we may write

$$x = \frac{1}{a_1 + x_1}$$

with $x_1 \in [0, 1]$. Then, since $x_1 \neq 0$, for otherwise x would be rational, we may define $a_2 = [1/x_1]$ and $x_2 = 1/x_1 - a_2$ to get

$$x = \frac{1}{a_1 + \frac{1}{a_2 + x_2}}$$

Inductively, we see that

$$x = \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\dots + \frac{1}{a_n + x_n}}}}$$

where $x_n = 1/x_{n-1} - a_n$ and $a_n = [1/x_{n-1}]$. This amounts to say that the sequence (x_n) is the trajectory of x under the *Gauss map* $G :]0, 1] \rightarrow]0, 1]$, defined as

$$x \mapsto 1/x - [1/x]$$

Observe that G is not defined at the origin, hence to iterate G we need to avoid all the preimages of 0, which are the rationals. This is not a problem if we want to study the statistical properties of G with respect to Lebesgue measure, since rationals form a subset of zero measure. The Gauss map admits an absolutely continuous invariant measure $\mu = \rho dx$, defined as

$$\mu(A) = \frac{1}{\log 2} \cdot \int_A \frac{1}{1+x} dx$$

for any Borel subset $A \subset]0, 1]$. The denominator $\log 2$ is there to normalize the measure, so we just have to check the invariance criterium for the density $\rho(x) = 1/(1+x)$. Since any $x' \in]0, 1]$ has one preimage $x_k = 1/(x' + k)$ in each interval $]1/(k+1), 1/k]$, we compute

$$\begin{aligned} \sum_{x \in G^{-1}\{x'\}} \frac{\rho(x)}{|G'(x)|} &= \sum_{k \geq 1} \frac{x_k^2}{1+x_k} \\ &= \sum_{k \geq 1} \left(\frac{1}{x' + k} - \frac{1}{x' + k + 1} \right) \\ &= \frac{1}{1+x'} = \rho(x') \end{aligned}$$

and we are done.

8 Recurrences

8.1 Limit sets and recurrent points

Omega and alpha limit sets. Let $f : X \rightarrow X$ be a continuous transformation of a topological space X . The simplest thing that an infinite (i.e. not periodic) orbit can do is to be (the image of a) convergent (trajectory). In this case, as we already know, the limit must be a fixed point of the transformation.

Trajectories, even when not convergent, may have at least convergent subsequences. The ω -limit set of a point $x \in X$ is the set $\omega_f(x) \subset X$ of those points $x' \in X$ such that there exists a sequence of times $n_i \rightarrow \infty$ (i.e. an increasing map $i \mapsto n_i$) such that $f^{n_i}(x) \rightarrow x'$ when $i \rightarrow \infty$. A little reflection shows that it is

$$\omega_f(x) = \bigcap_{n=0}^{\infty} \overline{\bigcup_{k \geq n} \{f^k(x)\}}$$

Observe that, if the orbit of x is not finite (i.e. if x is not periodic), then the ω -limit set of x is the derived set of its forward orbit, i.e. $\omega_f(x) = O_f^+(x)'$. It is clear that $\omega_f(x)$ is a closed (possibly empty) and $+$ -invariant subset of X .

$\text{Lim}(f) = \bigcup_{x \in X} \omega_f(x)$ denotes the set of ω -limit points of all the $x \in X$. If x is periodic, then $\omega_f(x)$ coincides with its orbit. There follows that

$$\text{Per}(f) \subset \text{Lim}(f).$$

If f is invertible, we may also define the α -limit set of $x \in X$ as $\alpha_f(x) := \omega_{f^{-1}}(x)$, i.e. the set of points $x' \in X$ such that there exists a sequence of times $n_i \rightarrow \infty$ such that $f^{-n_i}(x) \rightarrow x'$ when $i \rightarrow \infty$. In this case, both $\omega_f(x)$ and $\alpha_f(x)$ are closed and invariant subsets of X . $\text{Lim}(f^{-1}) = \bigcup_{x \in X} \alpha_f(x)$ denotes the set of all α -limit points of an invertible transformation f .

Limit sets in compact spaces. Both the ω and the α -limit sets of a generic point can be empty. For example, all the limit points for the translation $f(x) = x + 1$ of the real line are empty.

This may happen, of course, only if the phase space X is not compact. Indeed, if X is compact, then the trajectory of any point admits convergent subsequences (by sequential compactness, which holds for compact metric spaces), and therefore $\omega_f(x) \neq \emptyset$ for all $x \in X$. For the same reason, if f is a homeomorphism of a compact metric space, $\alpha_f(x) \neq \emptyset$ for all points $x \in X$. In particular, the sets $\text{Lim}(f^{\pm 1})$ are not empty.

ex: Show that $\omega_f(x)$ is closed and $+$ -invariant. Show that if f is a homeomorphism, then $\omega_f(x)$ and $\alpha_f(x)$ are closed and invariant.

ex: Give examples such that $\omega_f(x)$ and $\alpha_f(x)$ are empty.

ex: Show that $\text{Per}(f) \subset \text{Lim}(f)$.

Recurrent points. Let $f : X \rightarrow X$ be a topological dynamical system. The point $x \in X$ is *recurrent* if $x \in \omega_f(x)$. It is clear that this is equivalent to asking that given any neighbourhood B of x there exists a time $n \geq 1$ such that $f^n(x) \in B$. Indeed, choosing smaller neighbourhoods (if $f^n(x) \neq x$, so that x is not already periodic with period n), this also implies that the trajectory of x passes infinitely often in any such neighbourhood, i.e. $f^n(x) \in B$ for infinitely many times $n \geq 1$.

$\text{Rec}(f)$ denotes the set of recurrent points for f . A periodic point is obviously recurrent, therefore

$$\text{Per}(f) \subset \text{Rec}(f).$$

If f is a homeomorphism, it also makes sense to consider the set $\text{Rec}_{f^{-1}}$, the set of those points $x \in X$ such that $x \in \alpha_f(x)$.

ex: Define a partial order in X as follows: $x \prec x'$ if for any neighbourhood U of x and V of x' there exists a time $n \geq 1$ such that $f^n(U) \cap V \neq \emptyset$. Show that x is recurrent iff $x \prec x$.

ex: Show that $\text{Per}(f) \subset \text{Rec}(f)$.

ex: Give examples which show that both $\text{Rec}(f)$ and $\text{Rec}(f^{-1})$ may be empty.

Non-wandering set. The point x is *wandering* (the greek word for “wandering” was $\pi\lambda\alpha\nu\eta\tau\eta\varsigma$, i.e. *planet*) if it admits a neighborhood which is disjoint from all its iterates, i.e. if there exists an open set U containing x such that $U \cap f^n(U) = \emptyset$ for all times $n \geq 1$. Conversely, the point x is not wandering if for any neighbourhood U of x there exists a time $n \geq 1$ such that $f^n(U) \cap U \neq \emptyset$.

The *non-wandering set* $\text{NW}(f)$ is the set of those points x which are not wandering. This is the set where the interesting dynamics takes place, since other points are “forgotten” as time passes.

The non-wandering set is closed (the set of wandering points is open by definition, since any point in a sufficiently small open neighbourhood of a wandering point is itself wandering) and $+$ -invariant. It contains the ω -limit points of all points in X , as well as the recurring points. Thus, the inclusions are

$$\text{Per}(f) \subset \text{Lim}(f) \subset \text{NW}(f) \quad \text{and} \quad \text{Per}(f) \subset \text{Rec}(f) \subset \text{NW}(f)$$

If f is an homeomorphism, $\text{NW}(f)$, which is equal to $\text{NW}(f^{-1})$, is also invariant, and contains the ω - and α -limits of all points of X .

If X is compact, then $\text{NW}(f) \neq \emptyset$, since any point $x \in X$ have $\omega_f(x) \neq \emptyset$ and $\text{Lim}(f) \subset \text{NW}(f)$.

ex: Show that the non-wandering set of a homeomorphism is closed, invariant, and contains the ω and α -limit sets.

ex: Show that if f is a homeomorphism, then $\text{NW}(f) = \text{NW}(f^{-1})$.

ex: Give examples that show that $\text{NW}(f)$ may be empty.

ex: Show that $\text{Per}(f) \subset \text{Rec}_f \subset \text{NW}(f)$. Give example that show that these inclusions may be strict.

ex: Show that $\text{Per}(f) \subset \text{Rec}(f) \subset \text{NW}(f)$ and therefore $\overline{\text{Per}(f)} \subset \overline{\text{Rec}(f)} \subset \text{NW}(f)$. More difficult is to find example showing that these inclusion may be strict.

ex: Find the non-wandering sets of linear maps of the plane.

8.2 Dirichlet theorem on Diophantine approximation

Rotations of the circle and Dirichlet theorem on Diophantine approximation. Consider a rotation $R_\alpha : x + \mathbb{Z} \mapsto x + \alpha + \mathbb{Z}$ of the circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$. If α is rational, all points are trivially recurrent, being periodic. When α is irrational, recurrence of a point $x + \mathbb{Z}$ means that for any $\varepsilon > 0$ there exist an infinity of times $q \in \mathbb{N}$ such that that

$$d(x + \mathbb{Z}, x + q\alpha + \mathbb{Z}) < \varepsilon$$

or, equivalently, that for any $\varepsilon > 0$ there exist an infinity of rationals p/q such that

$$|q\alpha - p| < \varepsilon \quad \text{i.e.} \quad \left| \alpha - \frac{p}{q} \right| < \frac{\varepsilon}{q}$$

Indeed, much more is true, and is a consequence of the following classical result by Dirichlet³³ on Diophantine approximation (see [HW59]).

³³L.G.P. Dirichlet, Verallgemeinerung eines Satzes aus der Lehre von den Kettenbrüchen nebst einigen Anwendungen auf die Theorie der Zahlen, *S. B. Preuss. Akad. Wiss.* (1942), 93-95.

Theorem 8.1 (Dirichlet, 1842). *For any number θ and any positive integer $Q \in \mathbb{N}$ there exist $p \in \mathbb{Z}$ and $q \in \mathbb{Z} \setminus \{0\}$ such that*

$$|q\theta - p| < 1/Q \quad \text{and} \quad |q| \leq Q. \quad (8.1)$$

and, a fortiori,

$$|\theta - p/q| < 1/q^2. \quad (8.2)$$

Proof. Divide the unit interval $[0, 1)$ into the Q subintervals

$$[0, 1/Q), [1/Q, 2/Q), [2/Q, 3/Q), \dots, [(Q-1)/Q, 1)$$

of equal length $1/Q$, and consider the $Q+1$ points³⁴

$$\{0\}, \{\theta\}, \{2\theta\}, \dots, \{Q\theta\}$$

inside $[0, 1)$. By the box principle (which Dirichlet stated to prove this theorem!), at least two of those points, say $\{k\theta\}$ and $\{k'\theta\}$ with $k > k'$, belong to the same subinterval. Therefore, there exist integers a, a' such that $|k\theta - a - (k'\theta - a')| < 1/Q$. The theorem follows taking $q = k - k'$ and $p = a - a'$, and observing that $q \leq Q$. \square

For rational θ , there are only finitely many integers q and p satisfying the above inequalities (8.1). Indeed, if $\theta = a/b$ and $p/q \neq a/b$ (we may assume that both are reduced fractions), then

$$|q\theta - p| = \frac{|qa - pb|}{|b|} \geq \frac{1}{|b|}$$

(because the numerator is the absolute value of a non-zero integer) and therefore no fraction different from a/b may satisfy the inequalities (8.1) if Q is larger than $|b|$.

On the other hand, if θ is irrational and $p_1/q_1, p_2/q_2, \dots, p_n/q_n$ are any finite number of fractions satisfying (8.2), we may consider an integer Q larger than the inverse of

$$\varepsilon = \min_{1 \leq k \leq n} |q_k\theta - p_k| > 0$$

and produce, by theorem 8.1, one more fraction p/q satisfying (8.2). Thus,

Theorem 8.2 (Dirichlet, 1842). *For any irrational number θ there exist infinitely many reduced fractions p/q such that*

$$|\theta - p/q| < 1/q^2.$$

In particular, any point $x + \mathbb{Z}$ is recurrent for an irrational rotation of the circle.

8.3 Poincaré recurrence theorem

If f satisfies a condition (natural in physics) like “preserving a probability measure”, then there are a lot of recurrent points, actually almost any point is recurrent. If, moreover, the probability measure is diffuse, i.e. any non-empty open set has positive measure, then the set of recurrent points is also dense. These results, discovered by Henri Poincaré around 1890, motivated the modern theory of dynamical systems. They show how weak informations on the transformation (or the flow of a differential equation) may yield significative qualitative information about “almost all” orbits of the system. Here follow the precise statements, together with all the necessary technical details. If you don’t know the meaning of some words, like “measurable” or “Borel set”, don’t worry, just try to understand what’s going on. Poincaré himself didn’t know, yet!

You may look at this wonderful lecture on Poincaré recurrence theorem by Etienne Ghys in YouTube: <https://www.youtube.com/watch?v=21fHNMccrY8#t=1741>

³⁴As usual, $\{x\}$ denotes the “fractional part” of x , so that any real number may be written as a sum $x = [x] + \{x\}$ for some unique $[x] \in \mathbb{Z}$ and $\{x\} \in [0, 1)$.

Theorem 8.3 (Poincaré recurrence theorem). *Let $f : X \rightarrow X$ be an endomorphism of a probability space (X, \mathcal{E}, μ) , and let $A \in \mathcal{E}$. Then the set*

$$A_{\text{rec}} = \{x \in A \text{ t.q. } f^n(x) \in_{i.o.} A\}$$

of those points of A whose orbit passes through A infinitely often has total probability, namely $\mu(A_{\text{rec}}) = \mu(A)$

Proof. For $k \geq 1$, let

$$B_k = \{x \in A \text{ s.t. } f^n(x) \notin A \forall n \geq k\}$$

be the set of those points of A which never return in A after $n \geq k$ iterates. Observe that $B_k = A \cap (\cap_{n \geq k} f^{-n}(X \setminus A))$ and that $A_{\text{rec}} = A \setminus (\cup_{k \geq 1} B_k)$. In particular, this shows that A_{rec} is measurable. It is clear that $f^{-nk}(B_k) \cap B_k = \emptyset$ for any $n \geq 1$, since a point in the intersection would be a point $x \in B_k \subset A$ such that $f^{kn}(x) \in A$, and $kn \geq k$, contradicting the definition of B_k . For the same reason, $f^{-nk}(B_k) \cap f^{-mk}(B_k) = \emptyset$ for any $n > m \geq 0$. Therefore, the sets $f^{-nk}(B_k)$, for $n \in \mathbb{N}$, are pairwise disjoint. They also have all the same measure $\mu(f^{-nk}(B_k)) = \mu(B_k)$, because μ is invariant. This implies that $\mu(B_k) = 0$, because

$$\sum_{n \geq 1} \mu(B_k) = \sum_{n \geq 1} \mu(f^{-nk}(B_k)) = \mu(\cup_{n \geq 1} f^{-nk}(B_k)) \leq \mu(X) = 1.$$

There follows that $\mu(A_{\text{rec}}) = \mu(A)$. □

Now, let $f : X \rightarrow X$ be a continuous transformation of a metrizable topological space X , and let μ be an invariant Borel probability measure. Assume that (the topology of) X admits a countable basis $(U_i)_{i \in \mathbb{N}}$. We can apply the above theorem 8.3 to every open set U_i , and this easily implies that the set of recurrent points has full measure, i.e.

$$\mu(\text{Rec}(f)) = 1.$$

In particular, since any set of full measure is dense in the support of a Borel measure, we get the following general result.

Theorem 8.4 (topological Poincaré recurrence theorem). *Let $f : X \rightarrow X$ be a continuous transformation of a separable metrizable topological space X . The support of any invariant Borel probability measure μ is contained in the closure of the set of recurrent points, namely*

$$\text{supp}(\mu) \subset \overline{\text{Rec}(f)}.$$

In particular, if f admits an invariant measure μ which is diffuse (i.e. gives positive measure to any nonempty open set) then the set of recurrent points is dense in X , namely

$$\overline{\text{Rec}(f)} = X.$$

Observe that if f is a homeomorphism, then the same applies to $\text{Rec}(f^{-1})$, and the support of any invariant Borel probability measure is contained in the closure of $\text{Rec}(f) \cap \text{Rec}(f^{-1})$.

If you don't like the above proof, here is another, perhaps more elementary, of the last statement.

Proof. (of the last statement of theorem 8.4) Assume that the continuous map $f : X \rightarrow X$ preserves a diffuse Borel probability measure μ . For each $n \geq 1$, let

$$R_n := \{x \in X \text{ s.t. } \exists k \geq 1 \text{ s.t. } d(f^k(x), x) < 1/n\}$$

be the set of “ $1/n$ -recurrent” points. It is plain that $\text{Rec}(f) = \cap_{n=1}^{\infty} R_n$. The sets R_n are clearly open. To show that Rec_f is dense we must therefore show that each R_n is dense, since then the Baire theorem implies that also their countable intersection is dense. So, take any nonempty ball $B = B_r(p)$ with diameter $2r < 1/n$. Its inverse images $f^{-1}(B)$, $f^{-2}(B)$, $f^{-3}(B)$, \dots have all the

same measure by invariance, which is positive, i.e. $\mu(B) > 0$ (because the measure μ is diffuse). Since $\mu(X) = 1$, the $f^{-n}(B)$, for $n = 0, 1, 2, 3, \dots$, cannot be pairwise disjoint. There follows that there exist $k > 0$ and $n \geq 0$ such that $f^{-(n+k)}(B) \cap f^{-n}(B) \neq \emptyset$, and this implies that B contains a $1/n$ -recurrent point (for a point x in the intersection has both images $f^n(x)$ and $f^{n+k}(x) = f^k(f^n(x))$ in B , hence at distance $< 1/n$). Since B was arbitrary, this proves that each R_n is dense, and Baire theorem implies that $\text{Rec}(f)$ is dense too. \square

8.4 Transitivity and minimality

Transitive transformations. Let X be a complete and separable metric space. A continuous transformation $f : X \rightarrow X$ is (*topologically*) *+transitive* if it satisfies one of the following equivalent conditions:

- i) for any two not-empty open sets $U, V \subset X$ there exist a time $n \geq 0$ such that $f^n(U) \cap V \neq \emptyset$
- ii) there exists a point $x \in X$ such that $\omega_f(x) = X$
- iii) there exists a residual subset $R \subset X$ of points x such that $\omega_f(x) = X$

Proof. (of the equivalence) The implications $\text{iii}) \Rightarrow \text{ii}) \Rightarrow \text{i})$ are obvious, since, if $\omega_f(x) = X$, then the trajectory of x visits infinitely often all non-empty open subsets of X .

To show that $\text{i}) \Rightarrow \text{iii})$, the first observation is that condition i) amounts to say that, for all not-empty open set V , its orbit $\bigcup_{n \geq 0} f^{-n}(V)$ is dense, and, moreover, its orbits $\bigcup_{n \geq k} f^{-n}(V) = \bigcup_{n \geq 0} f^{-n}(f^{-k}(V))$ are also dense for all $k \geq 0$. Let $(U_i)_{i \in \mathbb{N}}$ be a countable basis for the topology of X . The family of $\bigcup_{n \geq k} f^{-n}(U_i)$, with $k \geq 0$ and $i \geq 1$, is a family of open and dense subsets of X . Its countable intersection $R = \bigcap_{i \in \mathbb{N}} \bigcap_{k \geq 0} \bigcup_{n \geq k} f^{-n}(U_i)$ is a residual set, and a point $x \in R$ has a trajectory which visits infinitely often all the open sets U_i , i.e. $\omega_f(x) = X$. \square

Also clear is that i) implies that X does not have isolated points (unless it has finite cardinality, trivial case in which X is a single orbit). This, in turns, implies that $O_f^+(x)' = X$ if $x \in R$.

ex: Prove the implications $\text{iii}) \Rightarrow \text{ii}) \Rightarrow \text{i})$ above.

ex: If $f : X \rightarrow X$ is *+transitive*, then $\text{NW}(f) = X$, since the non-wandering set contains the ω -limit sets of points $x \in X$.

ex: If $f : X \rightarrow X$ is *+transitive*, then $\text{Rec}(f)$ is a residual set (observe that if $\omega_f(x) = X$ then $x \in \omega_f(x)$).

Transitive homomorphisms. There exists a weaker notion, only meaningful for invertible transformations. A homeomorphism $f : X \rightarrow X$ is (*topologically*) *transitive* if it satisfies one of the following three conditions:

- i) for any two not-empty open sets $U, V \subset X$ there exists a time $n \in \mathbb{Z}$ such that $f^n(U) \cap V \neq \emptyset$
- ii) there exists a point $x \in X$ with dense orbit, i.e. such that $\overline{O_f(x)} = X$
- iii) there exists a residual set of points $x \in X$ with dense orbits, i.e. such that $\overline{O_f(x)} = X$

Proof. (of the equivalence) The implications $\text{iii}) \Rightarrow \text{ii}) \Rightarrow \text{i})$ are obvious, since if the full orbit $O_f(x)$ of x is dense, it visits at least once each not-empty open subset of X .

To show that $\text{i}) \Rightarrow \text{iii})$, we first observe that condition i) is equivalent to say that the orbit $\bigcup_{n \in \mathbb{Z}} f^n(V)$ of any not-empty open set V is dense in X . Let $(U_i)_{i \in \mathbb{N}}$ be a countable basis for the topology of X . The family $U_i^\pm = \bigcup_{n \in \mathbb{Z}} f^n(U_i)$ is therefore a family of dense and open sets. Its countable intersection $R = \bigcap_{i \in \mathbb{N}} U_i^\pm$ is a residual set, and the complete orbit of any point $x \in R$ visits at first once any of the open sets U_i . Therefore, if $x \in R$ then $\overline{O_f(x)} = X$. \square

Observe that $f : X \rightarrow X$ is a transitive homeomorphism iff f^{-1} is a transitive homeomorphism. A transitive homeomorphisms need not be also $+$ -transitive. Indeed, it may have no recurrent points and an empty non-wandering set, provided X is not compact!

It is interesting to observe that transitivity is a kind of “dynamical connectedness”, in the following precise sense.

Theorem 8.5. *A homeomorphism $f : X \rightarrow X$ is transitive iff X does not contain the disjoint union of two open invariant not-empty sets.*

Proof. The implication \Rightarrow is obvious. To show the reverse implication \Leftarrow , observe that, if $U, V \subset X$ are two not-empty open sets then $U^\pm = \bigcup_{n \in \mathbb{Z}} f^n(U)$ and $V^\pm = \bigcup_{n \in \mathbb{Z}} f^n(V)$ are open, not-empty invariant sets. If $U^\pm \cap V^\pm \neq \emptyset$, there exist times $n, m \in \mathbb{Z}$ such that $f^n(U) \cap f^m(V) \neq \emptyset$, which implies $f^{n-m}(U) \cap V \neq \emptyset$. \square

A consequence is the following useful criterium to decide when a homeomorphism cannot be transitive.

Theorem 8.6. *If $f : X \rightarrow X$ is a transitive homeomorphism, then all continuous invariant observable $\varphi : X \rightarrow \mathbb{R}$ is constant.*

Proof. Indeed, if φ is not constant, then it takes at least two values, say $a < b$. But the, if $c = (a + b)/2$, both $U = \{\varphi < c\}$ and $V = \{\varphi > c\}$ are invariant open disjoint and not-empty sets. \square

ex: Give examples of homeomorphisms $f : X \rightarrow X$ which are transitive but not $+$ -transitive.

ex: Show that a homeomorphism $f : X \rightarrow X$ is $+$ -transitive iff is transitive and its non-wandering set is the whole X (the implication \Leftarrow is obvious).

ex: It may happen that a transformation $f : X \rightarrow X$ is $+$ -transitive but some iterate f^n , with $n > 1$, is not. A trivial example is a permutation of a finite set, since some iterate is the identity transformation. In general, if X is compact, what happens is the following. There exist some finite covering $X = X_1 \cup X_2 \cup \dots \cup X_k$, where k divides n and the X_i 's are compact sets with no dense intersections $X_i \cap X_j$ if $i \neq j$, such that $f(X_i) = X_{i+1 \bmod k}$ and the restrictions $f^n|_{X_i}$ are $+$ -transitive. The idea is to choose a point $x \in X$ such that $\omega_f(x) = X$, and then define $X_i = \omega_{f^n}(f^i(x)) \dots$

Minimal sets. Let $f : X \rightarrow X$ be a continuous transformation. A closed not-empty $K \subset X$ is *minimal* if it is $+$ -invariant and if it does not contain proper closed $+$ -invariant subsets.

The orbit of a periodic point is an example of a minimal set.

If K is minimal, then the orbit of any $x \in K$ is dense in K , for otherwise its closure $\overline{\omega_f^+(x)}$ would be a proper $+$ -invariant closed subset of K . This implies that $x \in \omega_f(x)$, and therefore that all points of a minimal set are recurrent. If $\text{Min}(f)$ denotes the union of all minimal subsets of X , the inclusions are

$$\text{Per}(f) \subset \text{Min}(f) \subset \text{Rec}(f)$$

Of course, an arbitrary transformation $f : X \rightarrow X$ may not admit any minimal subsets. This is the case of a non-trivial translation $x \mapsto x + a$ of the real line.

If X is compact, we may consider the family \mathcal{C} of those subsets $C \subset X$ which are closed, not-empty and $+$ -invariant, equipped with the natural partial order given by inclusion “ \subset ”. The family is not empty, since it contains X itself. By Zorn lemma, \mathcal{C} contains a minimal element K , which is clearly a minimal set. More generally, we proved the following result.

Theorem 8.7. *If a continuous transformation $f : X \rightarrow X$ admits a compact $C \subset X$ such that $f(C) \subset C$, then it admits at least a minimal subset $K \subset C$.*

A consequence is that a continuous transformation $f : X \rightarrow X$ defined in a compact space X admits at least one recurrent point (which may be unique), since $\text{Min}(f) \subset \text{Rec}(f)$.

Minimal transformations. A continuous transformation $f : X \rightarrow X$ is *minimal* if it satisfies one of the following equivalent conditions:

- i) all forward orbits $O_f^+(x)$ are dense in X
- ii) X does not contain a proper closed and $+$ -invariant subset, and therefore is a minimal set.

Clearly, a minimal transformation is $+$ -transitive. If X is a discrete space, minimality implies that X is made of a single orbit, which may be finite. If X is not discrete, a minimal transformation cannot have periodic points.

Minimal homeomorphisms. A homeomorphism $f : X \rightarrow X$ is *minimal* if it satisfies one of the following equivalent conditions:

- i) all orbits $O_f(x)$ are dense in X
- ii) X does not contain a proper closed and invariant subset.

Minimal homeomorphisms are transitive. The discussion we made above about minimal sets may be repeated in this context. In particular, a homeomorphism $f : X \rightarrow X$ defined in a compact space X admits at least a minimal set $K \subset X$, which in this case is a closed not-empty invariant set which does not contain any proper closed invariant subsets.

For translations on a topological group, minimality coincides with topological transitivity.

Theorem 8.8. *A topological transitive translation on a topological group is minimal.*

Proof. Let G be a topological group, and consider the left translation $L_h : g \mapsto hg$ by an element $h \in G$. For any two elements $g, g' \in G$ and any time $n \in \mathbb{Z}$ we have $h^n g' = h^n g (g^{-1} g')$. Therefore, the orbit of g' is a right translation of the orbit of g , i.e. $O_{L_h}(g') = O_{L_h}(g) g^{-1} g'$. In particular, one orbit is dense iff all other orbits are dense. \square

ex: Give examples of transformations $f : X \rightarrow X$ such that $\text{Min}(f) = \emptyset$.

ex: Prove the implications i) \Leftrightarrow ii) above in the definition of “minimal transformation”.

ex: Prove the implications i) \Leftrightarrow ii) above in the definition of “minimal homeomorphism”.

8.5 Kronecker theorem on irrational rotations

Irrational rotations of the circle. A non-homogeneous version of Dirichlet’s theorem 8.1 was discovered by Kronecker. In its original formulation³⁵, it says that, given an irrational α , for any integer $Q > 0$ and any $y \in \mathbb{R}$ there exist integers p and $q > Q$ such that

$$|q\alpha - p - y| < 3/q \quad (8.3)$$

Let $R_\alpha(x + \mathbb{Z}) = x + \alpha + \mathbb{Z}$ denotes the rotation of the circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ by the irrational angle $\alpha \notin \mathbb{Q}$. If we don’t mind about the exact bound $3/q$ for the error, it says that for all $x + \mathbb{Z}$ and $x' + \mathbb{Z}$ in \mathbb{T} (the y above is $x' - x$) and any precision $\varepsilon > 0$ there exists a time $q > 0$ such that $d(R_\alpha^q(x + \mathbb{Z}), x' + \mathbb{Z}) < \varepsilon$. In our language,

Theorem 8.9 (Kronecker, 1884). *An irrational rotation of the circle is minimal (i.e. all its orbit are dense in the circle).*

Different proofs are presented in [HW59] (XXIII, Theorems 438, 439 and 440). Here we give just two.

Proof. Let $F \subset \mathbb{T}$ be the closure of an orbit of an irrational rotation of the circle. If F is not the whole circle, then its complement $I = \mathbb{T} \setminus F$, which is a not-empty open set, is a countable union of open intervals (arcs of the circle). Let J be (one of) the intervals of I of maximal length (why does it exist?), say $|J| > 0$. We claim that its images $f^n(J)$, with $n \in \mathbb{Z}$, are pairwise disjoint. Indeed, two such intervals $f^n(J)$ and $f^m(J)$, with $n \neq m$, cannot coincide, for otherwise

³⁵L. Kronecker, Die Periodensysteme von Funktionen Reeller Variablen, *Berliner Sitzungsberichte* (1884), 1071-1080.

the boundary points would be periodic points of the rotation (which is irrational), and also cannot have not-empty intersection, for otherwise their union would be an interval of I of bigger length. Since the rotation preserves the lengths, all $f^n(J)$ have the same positive length $|J|$, and this is impossible because the circle has finite (unit) length. \square

A less abstract proof (with a worse constant than in Kronecker's statement), along the ideas of Dirichlet theorem, is as follows.

Proof. Given $\varepsilon > 0$, chose a positive integer Q such that $1/Q < \varepsilon$. Divide the circle into Q intervals of length $1/Q$. By the box principle, at least two of the $Q + 1$ points

$$x + \mathbb{Z}, \quad R_\alpha(x + \mathbb{Z}), \quad \dots, \quad R_\alpha^q(x + \mathbb{Z})$$

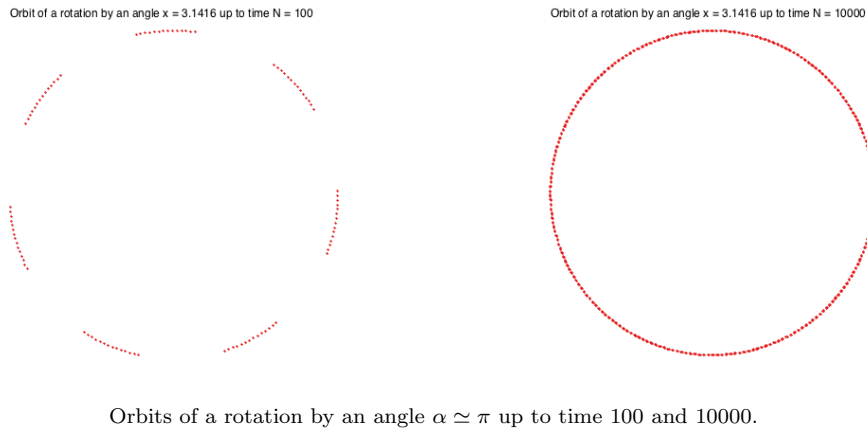
belong to the same interval, say $R_\alpha^i(x + \mathbb{Z})$ and $R_\alpha^j(x + \mathbb{Z})$ with $0 \leq i < j \leq Q$. Since rotations are isomerics,

$$d(R_\alpha^i(x + \mathbb{Z}), R_\alpha^j(x + \mathbb{Z})) = d(R_\alpha^k(x + \mathbb{Z}), x + \mathbb{Z}) < 1/Q < \varepsilon,$$

with $1 \leq k = j - i \leq Q$. Thus, the rotation R_α^k displaces points a (positive) distance $< \varepsilon$. It is clear then that the images $R_\alpha^{nk}(x + \mathbb{Z})$, with $n \geq 0$, pass infinitely often in a ε -neighbourhood of each point $x' + \mathbb{Z}$ of the circle. Thus, there exist integers $q = nk > Q$ and p such that $|q\alpha - p - y| < 1/Q < \varepsilon$, where $y = x' - x$. \square

As explained in [HW59], Kronecker theorem has a nice physical interpretation. It implies that orbits in a square billiard are either periodic, if the angle of incidence of the first hit to the boundary is a rational multiple of π , or dense in the square, otherwise. This is just the starting point of modern theory of billiards, a major area in dynamical systems.

The theorem has also an “algebraic” side. Observe that the orbit of $0 + \mathbb{Z}$, the identity of the abelian group \mathbb{R}/\mathbb{Z} , is the cyclic subgroup generated by $\alpha + \mathbb{Z}$. Therefore Kronecker theorem says that *the closed and proper subgroups of \mathbb{R}/\mathbb{Z} are the finite subgroups*.



Example of a non-measurable set. If you believe the axiom of choice, you may consider a set B made of one (exactly one!) point for any orbit of an irrational rotation R_α of the circle. The images $B_n = R_\alpha^n(B)$, for $n \in \mathbb{Z}$, are pairwise disjoint and cover the circle. If B , hence all its images, were Lebesgue-measurable, then

$$\sum_{n \in \mathbb{Z}} |B| = \sum_{n \in \mathbb{Z}} |B_n| = |\cup_{n \in \mathbb{Z}} B_n| = |\mathbb{R}/\mathbb{Z}| = 1$$

since rotations preserve Lebesgue measure, so that $|B_n| = |B|$. But there exists no size $b = |B| \geq 0$ such that $\sum_{n \in \mathbb{Z}} b = 1$.

ex: Also instructive is to see why rational rotations are not transitive, using theorem 8.6, since this extends to the higher dimensional torus. If $\alpha = p/q$ is rational, then the function $\varphi(x + \mathbb{Z}) = \sin(2\pi qx)$ is well defined in the circle \mathbb{R}/\mathbb{Z} , continuous, non-constant, and clearly invariant under the rotation R_α .

Rotations of the torus. Kronecker's theorem is actually much more general [?]. We say that the frequencies/numbers $\omega_1, \omega_2, \dots, \omega_k$ are *linearly independent over the rationals* if the only rational solution of the equation

$$n_1\omega_1 + n_2\omega_2 + \dots + n_k\omega_k = 0$$

is the trivial solution $n_1 = n_2 = \dots = n_k = 0$. An important example: the logarithms $\omega_k = \log p_k$ of different primes p_k are linearly independent, as follows from the uniqueness of prime decomposition.

Theorem 8.10 (Kronecker, 1884). *Let $\theta = (\theta_1, \theta_2, \dots, \theta_N) \in \mathbb{R}^N$. If $\theta_1, \theta_2, \dots, \theta_N, 1$ are linearly independent over the rationals, then any orbit*

$$x + \mathbb{Z}\theta + \mathbb{Z}^N$$

with $x \in \mathbb{R}^N$ is dense in the torus $\mathbb{T}^N = \mathbb{R}^N / \mathbb{Z}^N$.

This means that for any $x = (x_1, x_2, \dots, x_N) \in [0, 1)^N \approx \mathbb{T}^N := \mathbb{R}^N / \mathbb{Z}^N$ and any precision $\varepsilon > 0$ we can find integers p_1, p_2, \dots, p_N and $q \in \mathbb{Z}$ such that $|q\theta_k - p_k - x_k| < \varepsilon$ for all $k = 1, 2, \dots, N$. Chapter XXIII of [HW59] contains some different proofs.

8.6 Circle homeomorphisms

While studying vector fields on the 2-dimensional torus $\mathbb{R}^2 / \mathbb{Z}^2$, Henri Poincaré³⁶ was led to the necessity to classify possible dynamics of circle homeomorphisms. A model is that of rotations, where the dichotomy between closed or dense orbits reflects the rationality of the “rotation angle”. He discovered an invariant which plays a similar role for generic homeomorphisms.

Homeomorphisms of the circle. A homeomorphism $F : \mathbb{R} \rightarrow \mathbb{R}$ such that $F(x+1) = F(x)+1$ defines a orientation preserving homeomorphism $f : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$, according to $f(x+\mathbb{Z}) := F(x)+\mathbb{Z}$. Conversely, a bijection $f : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ is an *orientation preserving homeomorphism* of the circle if there exists a homeomorphism $F : \mathbb{R} \rightarrow \mathbb{R}$ such that $F(x+1) = F(x)+1$ (observe that this condition implies that F is strictly increasing) and $f(x+\mathbb{Z}) = F(x)+\mathbb{Z}$. Such F is called *lift* of f . Clearly, the lift is not unique, but any two lifts F and G of f differs by an integer constant, i.e. $F(x) = G(x) + n$ for some $n \in \mathbb{Z}$.

For example, a lift of the rotation $R_\alpha(x) = x + \alpha + \mathbb{Z}$ is the translation $F(x) = x + \alpha$. It is clear that if ε is sufficiently small, then $F_\varepsilon(x) = x + \alpha + \varepsilon \sin(2\pi x)$ induces a homeomorphism of the circle $f_\varepsilon(x+\mathbb{Z}) = x + \alpha + \varepsilon \sin(2\pi x) + \mathbb{Z}$, that may be considered a small variation of the rigid rotation $f_0 = R_\alpha$.

Rotation number. Let $f : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ be a orientation preserving homeomorphism of the circle, and let $F : \mathbb{R} \rightarrow \mathbb{R}$ be one of its lifts.

The *rotation number* of f is the “angle”

$$\rho(f) := \tau(F) + \mathbb{Z} \in \mathbb{R}/\mathbb{Z} \quad (8.4)$$

where $\tau(F)$ is the *translation number* of F , defined by

$$\tau(F) := \lim_{n \rightarrow \infty} \frac{F^n(x) - x}{n} \quad (8.5)$$

where x is an arbitrary point of the line. This makes sense once we prove that the limit exists and does not depend on the initial point x , and that its class in the circle does not depend on the particular lift.

For example, the translation number of the translation $F(x) = x + \alpha$ is α , and therefore the rotation number of the rotation R_α is $\alpha + \mathbb{Z}$.

The main ingredient of the existence proof is the following fact, of independent importance, about subadditive sequences.

³⁶H. Poincaré, Sur les courbes définies par les équations différentielles, *J. Math. Pures App.* Série IV **1** (1885), 167-244.

Theorem 8.11 (subadditive sequence lemma). *Let $(a_n)_{n \in \mathbb{N}}$ be a quasi-subadditive real sequence, i.e. a sequence such that*

$$a_{n+m} \leq a_n + a_m + c$$

for any $n, m \in \mathbb{N}$ and some $c \geq 0$. Then there exists the limit

$$\lim_{n \rightarrow \infty} \frac{a_n}{n} \in \mathbb{R} \cup \{-\infty\}.$$

Proof. It is clear that existence of the limit $\lim_{n \rightarrow \infty} a_n/n$ is equivalent to existence of the limit $\lim_{n \rightarrow \infty} b_n/n$ for the sequence $b_n = a_n + c$. The sequence (b_n) is now subadditive, i.e.

$$b_{n+m} \leq b_n + b_m.$$

Subadditivity implies $b_n \leq nb_1$. Thus the sequence (b_n/n) is bounded above, hence there exists $\lambda = \liminf_{n \rightarrow \infty} b_n/n < \infty$. Given $\varepsilon > 0$, there exists a natural m such that $b_m/m < \lambda + \varepsilon$. A generic positive integer as $n = km + r$, with k a non-negative integer and $0 \leq r < m$. Let $B = \max_{1 \leq i < m} b_i$. Subadditivity also implies

$$\begin{aligned} b_n/n &\leq (b_{km} + b_r)/n \leq (kb_m + b_r)/n \\ &\leq b_m/m + b_r/n \leq \lambda + \varepsilon + B/n \end{aligned}$$

By the arbitrariness of ε , the inequality above implies that $\limsup_{n \rightarrow \infty} b_n/n \leq \lambda$. Thus, the limit $\lim_{n \rightarrow \infty} b_n/n$ exists and is equal to λ . \square

Theorem 8.12. *The limit $\tau(F)$ in (8.5) exists.*

Proof. The lift F and its iterates F^n are increasing homeomorphisms of the real line satisfying $F^n(x+1) = F^n(x) + 1$ for all $x \in \mathbb{R}$. In particular, $F^n(x) - x$ are periodic functions of period one. This implies that

$$\max_{x, x'} |(F^n(x) - x) - (F^n(x') - x')| \leq 1$$

since, by periodicity, we may compute the maximum inside the unit interval $[0, 1]$, and we know that F^n is increasing and that the image $F^n([0, 1])$ is an interval of unit length. Let $a_n = F^n(x) - x$. The above inequality implies that the sequence (a_n) is quasi-subadditive, i.e.

$$a_{n+m} \leq a_n + a_m + c$$

for all $n, m \geq 0$ and some constant c . Indeed,

$$\begin{aligned} F^{n+m}(x) - x &= F^n(F^m(x)) - F^m(x) + F^m(x) - x \\ &= F^n(x) - x - F^n(x) + x + F^n(F^m(x)) - F^m(x) + F^m(x) - x \\ &\leq F^n(x) - x + F^m(x) - x + 1 \end{aligned}$$

so that we may choose $c = 1$. The theorem follows from theorem 8.11. \square

Theorem 8.13. *The limit $\tau(F)$ in (8.5) does not depend on x .*

Proof. We already saw that $|(F^n(x) - x) - (F^n(x') - x')| \leq 1$. Therefore,

$$\left| \frac{F^n(x) - x}{n} - \frac{F^n(x') - x'}{n} \right| \leq 1/n$$

for all x, x' and all n . This implies that $\tau(F)$ is independent on the chosen point x . \square

Theorem 8.14. *The class $\rho(f) = \tau(F) + \mathbb{Z}$ does not depend on the lift F of f .*

Proof. Two different lifts, say F and G , differ by an integer, i.e. $G(x) = F(x) + k$ for some $k \in \mathbb{Z}$. This implies that $G^n(x) - x = F^n(x) - x + nk$ and therefore that $\tau(F) = \tau(G) + k$. \square

Finally,

Theorem 8.15. *The rotation number $\rho(f)$ is invariant under topological conjugations.*

Proof. Let $h : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ be a conjugation between the homeomorphisms f and g . If H is a lift of h and F is a lift of F , then $H \circ F \circ H^{-1}$ is a lift of g . But the difference $(H \circ F \circ H^{-1})^n(x) - F^n(x)$ is bounded, independently of x and n . Indeed, we may observe that

$$|(H \circ F \circ H^{-1})^n| = |H \circ F^n \circ H^{-1}|, \quad |H(x) - x| \quad \text{and} \quad |H^{-1}(x) - x|$$

are bounded by a constant which does not depend on x , and use triangular inequality. This implies that $\tau(F) = \tau(H \circ F \circ H^{-1})$, and therefore that $\rho(f) = \rho(g)$. \square

Of course, the rotation number of a rotation $R_\alpha(x + \mathbb{Z}) = x + \alpha + \mathbb{Z}$ is α itself.

ex: Show that $\rho(f^q) = q \cdot \rho(f) + \mathbb{Z}$ (observe that, if F is a lift of f , then F^n is a lift of f^n).

Poincaré classification theorem The rotation number contains the following information about the dynamics of an homeomorphism.

Theorem 8.16 (Poincaré). *The rotation number $\rho(f)$ is rational iff the homeomorphism f admits periodic points.*

Proof. (\Leftarrow) If $F^q(x) = x + p$ with integers $q \geq 1$ and p , then $F^{nq}(x) - x = np$ for all n , and therefore $\tau(F) = p/q$.

(\Rightarrow) Observing that $\rho(f^q) = q \cdot \rho(f) \bmod \mathbb{Z}$, it is sufficient to prove that $\rho(f) = 0$ implies that f has a fixed point. Now, if f does not have fixed points and F is a lift of f , then the function $F(x) - x$ has values in \mathbb{R}/\mathbb{Z} . But the image of the real line by a continuous function is an interval. Therefore, there exists a lift F such that $F(x) - x$ takes values in the open unit interval $(0, 1)$. Since $F - \text{id}$ is periodic with period one, its maximum and minimum are both different from 0 and 1. Thus, there exists $\varepsilon > 0$ such that $\varepsilon < F(x) - x < 1 - \varepsilon$ for all $x \in [0, 1]$. Iterating, this implies $n\varepsilon < F^n(0) < n(1 - \varepsilon)$ and therefore that $\tau(F)$ is not integer. \square

Indeed, one can also prove that if $\rho(f)$ is rational then all periodic points share the same period. Thus, in order to understand the structure of orbits of a homeomorphism with rational rotation number is sufficient to study the case of zero rotation number, i.e. homeomorphisms with a fixed point. If $C = \text{Fix}(f)$ is the set of fixed points (which may be any compact subset of the circle), then f induces homeomorphism in any connected component I of the open set $\mathbb{R}/\mathbb{Z} \setminus C$. Images $f^n(x)$ of points $x \in I \subset (\mathbb{R}/\mathbb{Z}) \setminus C$ converge to points in $\partial I \subset C$ when $n \rightarrow \pm\infty$.

The dynamics of homeomorphisms with irrational rotation number is described by the following result.

Theorem 8.17 (Poincaré). *Let $f : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ be a orientation preserving homeomorphism with irrational rotation number. Then*

- i) *either f is minimal, i.e. the orbit of all points are dense in the circle,*
- ii) *or there exist an invariant compact subset $K \subset \mathbb{R}/\mathbb{Z}$, perfect and with empty interior (i.e. a Cantor set) such that the ω -limitset of all points of the circle is equal to K .*

Proof. By Zorn lemma, the family of not-empty compact invariant sets of the circle, equipped with the natural partial order given by inclusion, admits a minimal set K . By minimality, the orbit of any point of K is dense in K . The boundary ∂K and the derived set K' are also compact and invariant, so they must be empty or coincide with K itself. Since f has no fixed points, K cannot be finite. By the Bolzano-Weierstrass theorem $K' \neq \emptyset$, hence $K' = K$, i.e. K is perfect.

Now, if $\partial K = \emptyset$, then $K = \mathbb{R}/\mathbb{Z}$ and therefore f is minimal. If, otherwise, $\partial K = K$, then K has empty interior. Let $x \in (\mathbb{R}/\mathbb{Z}) \setminus K$ let I be the connected component of $(\mathbb{R}/\mathbb{Z}) \setminus K$ which contains x . The images $f^n(I)$ are pairwise disjoint (because f has no fixed points), and therefore $\text{diam}(f^n(I)) \rightarrow 0$ when $n \rightarrow \infty$. If $x' \in \partial I \subset K$, then $\omega_f(x') = K$, and the previous observation implies that also $\omega_f(x) = K$, because $d(f^n(x), f^n(x')) \leq \text{diam}(f^n(I)) \rightarrow 0$ when $n \rightarrow \infty$. In particular, this shows that the minimal set K is unique. \square

More interesting is the following result, also due to Poincaré.

Theorem 8.18 (Poincaré classification theorem). *Let $f : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ be a orientation preserving homeomorphism with irrational rotation number ρ .*

- i) If f is minimal, then it is topologically conjugated to the rotation R_ρ .*
- ii) If f is not minimal, then the rotation R_ρ is a factor of f .*

Indeed, if f is minimal we may construct a conjugation H between one orbit of f and one orbit of R_ρ , and then define the full conjugation h by continuity, using the fact that orbits are dense. This is possible because orbits of f have the “same order” than orbits of a rotation. If f is not minimal, it is possible to construct a semiconjugation $h : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ such that \mathbb{R}/\mathbb{Z} is the image $h(K)$ of the minimal set of f . Somehow, the semiconjugation “forgets” $(\mathbb{R}/\mathbb{Z}) \setminus K$, the wandering set of f .

This was a starting point of a beautiful story, starting with Denjoy in the 30's of the last century, and due to mathematicians like Michaël Hermann, Adrien Douady, Jean-Christophe Yoccoz, ... See [MS93].

9 Chaos

9.1 Sensitive dependence on initial conditions

Regular points and loss of memory. Iterations of a continuous transformation $f : X \rightarrow X$ of a metric space divide in a natural manner the phase space into two classes of points, depending whether orbits are stable or unstable under small perturbations.

The point $x \in X$ is *regular* if the family $\{f^n\}_{n \geq 0}$ is equicontinuous at x , i.e. if for all $\varepsilon > 0$ there is a neighbourhood B of x such that for all $x' \in B$ and all times $n \geq 0$

$$d(f^n(x), f^n(x')) < \varepsilon$$

So, one orbit for each ε -ball contain informations on orbits of all regular points. In particular, if X is compact and all point is regular, we only need a finite number of orbits to describe the dynamics up to an error ε .

The point $x \in X$ is *not regular* if there exists $\delta > 0$ such that in any neighbourhood B of x there are points $x' \in B$ such that

$$d(f^n(x), f^n(x')) > \delta$$

for some time $n \geq 1$. This means that f has “sensitive dependence on initial conditions” near x . In some sense, trajectories of points nearby x “lose their memory” of x . If the set of non-regular points is compact, the δ above may be chose uniformly for all of the points. This suggest the following definition.

Equicontinuous homeomorphisms. Let $f : X \rightarrow X$ be a equicontinuous homeomorphism of a compact metric space (X, d) , so that for all $\varepsilon > 0$ there exists a $\delta > 0$ such that $d(x, y) < \delta$ implies $d(f^n(x), f^n(y)) < \varepsilon$ for all times $n \in \mathbb{Z}$. Define

$$d^\infty(x, y) := \sup_{n \in \mathbb{Z}} d(f^n(x), f^n(y)).$$

It is clear that d^∞ is a metric on X , and that f is an isometry of (X, d^∞) . Thus, equicontinuous homeomorphisms of compact metric spaces behave as isometries.

Sensitive dependence on initial conditions. The continuous transformation $f : X \rightarrow X$ has *sensitive dependence on initial conditions* if all points are uniformly not-regular, i.e. if there exists $\delta > 0$ such that for all $x \in X$ and all neighbourhoods B of x , there exist $x' \in B$ and a time $n \geq 1$ such that

$$d(f^n(x'), f^n(x)) > \delta$$

In other words, no matter how small our sensibility ε is, in a ε -neighbourhood of any point x there is another point x' such that the futures of x and x' is uncorrelated, being at macroscopic (relative to ε) distance δ after some time n . Thus, a small change in the initial conditions may produce a large change at later times, a phenomenon popularized as “butterfly effect” by Edward Lorenz.

Of course this phenomenon is unexpected when the phase space is compact, for otherwise there is plenty of space for orbits to diverge from each other. Observe also that sensitive dependence is not compatible with preserving distances, hence isometries (like rotations of a torus) cannot have such a property. Thus, in order to display this kind of chaotic behaviour, a map must somehow “stretch” and “fold”, as our examples below will show.

Chaos. The combination of sensitive dependence on initial condition and a dense set of periodic points is usually referred as *chaos*³⁷.

³⁷The Greek word $\chi\alpha\omicron\varsigma$, which we may translate as “abyss”, contains the same root $\chi\alpha$ - (and probably comes from) of the verbs $\chi\alpha\lambda\upsilon\epsilon\iota\nu$ and $\chi\alpha\sigma\chi\epsilon\iota\nu$, which mean “open-itself”, “open the mouth” or “yawn” (cfr. $\chi\alpha\sigma\mu\alpha$, i.e. “chasm”). It has been used in some greek cosmogonies to mean “the desordered mixture of elements preceeding the formation of the $\chi\omicron\sigma\mu\omicron\sigma$, the ordered universe”.

Julia and Fatou sets. The above dichotomy is particularly meaningful for endomorphisms of the Riemann sphere $\overline{\mathbb{C}} = \mathbb{C} \cup \{\infty\}$, rational transformations $f : \overline{\mathbb{C}} \rightarrow \overline{\mathbb{C}}$

$$z \mapsto f(z) = \frac{p(z)}{q(z)}$$

where p and q are polynomials. A point $z \in \overline{\mathbb{C}}$ is *regular* if it admits a neighbourhood U such that $\{f^n|_U\}_{n \geq 1}$ is a “normal family” (i.e. every sequence of elements of the family admits a locally uniformly convergent subsequence).

The set $F \subset \overline{\mathbb{C}}$ of regular points, which is an open subset of the Riemann sphere, is called *Fatou set*. The complementary set, the compact $J = \overline{\mathbb{C}} \setminus F$, is called *Julia set*. The Julia set is where the interesting, i.e. disordered, dynamics takes place.

For example, if $f(z) = z^n$, then the Julia set is the unit circle \mathbb{S} . For perturbations, like for example $f(z) = z^2 + c$, the Julia set becomes a very irregular curve, typically of Hausdorff dimension > 1 , or a “dust” like a Cantor set.

The investigation on the dynamics of rational maps started at the beginning of the last century (1918-19) with Gaston Julia and Pierre Fatou. The contemporary theory, due essentially to sophisticated ideas in complex analysis, is one of the greatest successes of the modern theory of dynamical systems, thanks to the work of Douady, Sullivan, Milnor, Yoccoz, Lyubich, McMulle, ... and many others. A great introduction is in a famous lectures notes by John Milnor [Mi91].

Lyapunov exponents. Sensitive dependence on initial conditions can be quantified when the phase space has a differentiable structure. Assume, to begin with, that $f : X \rightarrow X$ is a smooth map of an interval $X \subset \mathbb{R}$. A small perturbation $x_0 + \delta x_0$ in the initial condition x_0 is amplified, under iterations of the map f , according to

$$\delta x_n = f^n(x_0 + \delta x_0) - f^n(x_0) \simeq (f^n)'(x_0) \cdot \delta x_0$$

up to first order. By the chain rule, the above derivative is a product

$$(f^n)'(x_0) = f'(x_{n-1}) \cdots f'(x_1) f'(x_0)$$

of the derivatives of f computed at the n -orbit of x_0 . If we expect an exponential growth of perturbations like

$$|\delta x_n| \simeq e^{\lambda n} |\delta x_0|,$$

we may estimate the *exponential growth rate* of the perturbations as $\lambda \simeq \frac{1}{n} \log |(f^n)'(x_0)|$ for large times n . Better, we may consider the limit as $n \rightarrow \infty$, and define the *Lyapunov exponent* of the (one-dimensional) map f at (the orbit of the point) x_0 as

$$\lambda(f, x_0) := \lim_{n \rightarrow \infty} \frac{1}{n} (\log |f'(x_{n-1})| + \cdots + \log |f'(x_1)| + \log |f'(x_0)|) \quad (9.1)$$

Of course, the above limit may not exist, but then we could consider the lim sup or the lim inf. On the other hand, the limit certainly exists for μ -almost all points, with respect to an invariant probability measure μ , as follows from the Birkhoff-Khinchin ergodic theorem 7.9 (provided the logarithm of the derivative of f is integrable). Clearly, a positive Lyapunov exponent implies sensitive dependence near the given initial condition.

The limit (9.1) certainly exists at periodic points. Indeed, if $f^n(x_0) = x_0$, then it is easy to see that the Lyapunov exponent exists and is equal to the arithmetic average of the logarithms of the absolute values of the derivatives of f at the finite orbit.

$$\lambda(f, x_0) = \frac{1}{n} (\log |f'(x_{n-1})| + \cdots + \log |f'(x_1)| + \log |f'(x_0)|)$$

In particular, it is negative if the orbit is attracting.

It is also clear that if the two maps $f : X \rightarrow X$ and $g : Y \rightarrow Y$ are conjugated by a smooth conjugation $\varphi : X \rightarrow Y$, i.e. $g = \varphi \circ f \circ \varphi^{-1}$, then corresponding points share the same Lyapunov exponent. Indeed, if $y_0 = \varphi(x_0)$, then $(g^n)'(y_0) = \varphi'(x_{n-1}) (f^n)'(x_0) (\varphi^{-1})'(y_0)$, and the

derivatives of φ disappear, after taking logarithms and dividing by n , when $n \rightarrow \infty$. Therefore, $\lambda(f, x_0) = \lambda(g, y_0)$, provided one of the two limits exists.

In dimension higher than one things get much more complicated. Indeed, if $f : X \rightarrow X$ is a smooth map defined in a open set $X \subset \mathbb{R}^N$ (or a differentiable manifold), then its derivative at a point $x \in X$ is a linear map $Df(x)$ from the tangent space $T_x X$ to the tangent space $T_{f(x)} X$, i.e. a $n \times n$ real matrix once fixed a coordinate system, possibly invertible (thus belonging to the general linear group $GL_N(\mathbb{R})$). The chain rule continues true in the form of the “cocycle identity”

$$Df^n(x_0) = Df(x_{n-1}) \dots Df(x_1) Df(x_0)$$

where products are rows by columns products (compositions of linear maps). But then, one could have different exponential rates along different directions, and therefore one should talk about a “Lyapunov spectrum”. Its existence is dealt by the *Oseledets’ multiplicative ergodic theorem*, and it is an interesting aspect of the contemporary research in dynamical systems, thanks to the work of Kaimanovich, Ruelle, Margulis, Katok, ...

ex: Consider the logistic map for different values of the parameter, and do simulations trying to estimate the Lyapunov exponents of some initial point.

9.2 Topological mixing

Topological mixing. A continuous map $f : X \rightarrow X$ is *topologically mixing* if for any two not-empty open sets $U, V \subset X$ there exists a time $n \geq 0$ such that for all times $k \geq n$

$$f^k(U) \cap V \neq \emptyset$$

This definition captures the idea that the future $f^k(U)$, with $k \geq n$, of any open set is “asymptotically independent” on its present, since it intersects stably all other not-empty open set V .

It is clear that a topologically mixing map is also + transitive. In particular, if f is topologically mixing then $NW_f = X$ and $\omega_f(x) = X$ is a generic property.

Theorem 9.1. *A topologically mixing map of a non trivial metric space has sensitive dependence on initial conditions.*

Proof. Indeed, let U and V be two disjoint open sets at distance at least $2\delta > 0$ (which exist if X contains at least two distinct points). Given $x \in X$, the orbit of any neighbourhood B of x intersects any not-empty open set starting from some time $n \geq 0$. This easily implies, by the triangular inequality, that there exists a point $x' \in B$ such that $d(f^n(x'), f^n(x)) > \delta$. \square

In particular, an isometry (as a torus rotation) cannot be topologically mixing.

ex: Does there exist a minimal (hence topologically transitive) homeomorphism which is not topologically mixing?

ex: Does there exist a topologically transitive map which is not minimal neither topologically mixing?

ex: (difficult) A continuous map $f : X \rightarrow X$ is *weakly mixing* if the product map $f \times f : X \times X \rightarrow X \times X$, defined by

$$(x, x') \mapsto (f(x), f(x'))$$

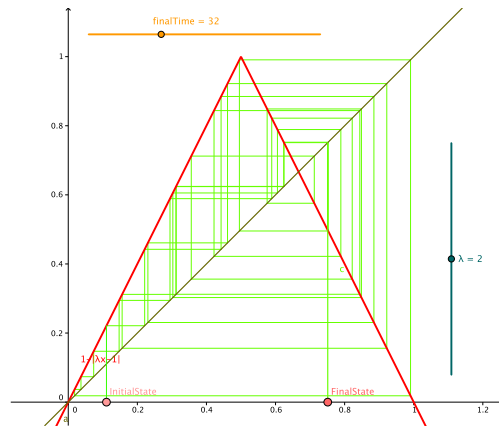
is topologically mixing. Show that a weakly mixing map of a non-trivial space X (i.e. which contains at least two points) has sensitive dependence on initial conditions. Show that all iterates f^n of a weakly mixing map of a compact space are +transitive. Show that

$$\text{mixing} \Rightarrow \text{weak mixing} \Rightarrow + \text{transitive}$$

and give examples which show that the reverse implications are false.

Tent map. One of the paradigms of chaotic maps is the *tent map* $T : [0, 1] \rightarrow [0, 1]$, defined by

$$T(x) = \begin{cases} 2x & \text{if } x < 1/2 \\ 2(1-x) & \text{if } x \geq 1/2 \end{cases}$$



Cobweb diagram of the tent map.

Iterations of the tent map are simple, since T is piecewise affine, and compositions of affine maps are affine maps. Indeed, it is clear (and not difficult to prove by induction) that in any dyadic interval as $I_{k,n} = [\frac{k}{2^n}, \frac{k+1}{2^n}]$, with $k = 0, 1, 2, \dots, 2^n - 1$, the iterate T^n is given by

$$T^n(x) = \begin{cases} 2^n \left(x - \frac{k}{2^n} \right) & \text{if } k \text{ is even} \\ 2^n \left(\frac{k+1}{2^n} - x \right) & \text{if } k \text{ is odd} \end{cases}$$

In particular, T^n is a strictly increasing or decreasing bijection of $I_{k,n}$ onto $[0, 1]$. The fixed point theorem then implies that T^n has exactly one fixed point in each of these intervals $I_{k,n}$ (which is repelling since the modulus of the derivative of T^n is $2^n > 1$, and only coincides with one of the boundary points when $k = 0$), and therefore the cardinality of n -periodic points is $P_n(T) = \text{card}(\text{Per}_n(T)) = 2^n$.

Moreover, since any not-empty open interval $U \subset [0, 1]$ contains one of the dyadic intervals $I_{k,n}$, if n is sufficiently large, and any such interval $I_{k,n}$ contains a periodic point of period d dividing n , periodic points are dense in the unit interval.

Finally, the tent map is topologically mixing. Indeed, since any not-empty open set $U \subset [0, 1]$ contains one of the $I_{k,n}$, its image under T^n is $T^n(U) = [0, 1]$, and, a fortiori, $T^m(U) = [0, 1]$ for all times $m \geq n$ because T is onto. This implies that $T^m(U) \cap V \neq \emptyset$ for all times $m \geq n$ and any other not-empty open set $V \subset [0, 1]$. Thus, there exists a residual set of points x such that $\omega_f(x) = [0, 1]$, i.e. with essentially unpredictable trajectory.

ex: Show that the tent map has Lyapunov exponent $\lambda = \log 2$ at almost all points.

ex: Show that $h : x \mapsto \sin^2(\pi x/2)$ is a topological conjugation between the tent map T and the transformation $f_4 : [0, 1] \rightarrow [0, 1]$ of the quadratic family, defined by $f_4(x) = 4x(1-x)$. This shows that f_4 has the same properties than T , e.g. it is topologically mixing and has a dense set of periodic points. Moreover, and this is surprising, this also provides explicit formulae for the trajectories of f_4 . Indeed, if the initial condition is $x_0 = \sin^2(\pi t_0)$, then $x_n = f_4^n(x_0)$ is given by $x_n = \sin^2(2^n \pi t_0)$.

ex: Discuss also the dynamics of the map $S : [0, 1] \rightarrow [0, 1]$ defined by

$$S(x) = \begin{cases} 2x & \text{if } x < 1/2 \\ 2x - 1 & \text{if } x \geq 1/2 \end{cases}$$

Observe that S is not continuous, but it is not much different from the tent map.

9.3 Expanding maps of the circle

The obvious way to force sensitive dependence on initial conditions is “stretching”, for example dilating distances at least along some directions, and “folding”, for example, taking quotients.

Expanding maps. A continuous transformation $f : X \rightarrow X$ of a metric space is *expanding* if there exist $\lambda > 1$ and $\varepsilon > 0$ such that for all $x, x' \in X$ at distance $d(x, x') < \varepsilon$ we have

$$d(f(x), f(x')) \geq \lambda \cdot d(x, x')$$

This condition, which looks like an opposite of a contraction, is a local condition, since otherwise compact spaces would not admit expanding maps with large ε . On the other side, it is precisely in compact phase spaces X that stretching induces chaotic orbits, since there is not enough space to escape, and divergent trajectories are forced to come back, eventually.

The mere existence of expanding maps also implies strong topological restrictions on the possible phase spaces. If X is a manifold, then its universal cover must be \mathbb{R}^N , and even then, its fundamental group cannot be arbitrary. For example, between all the orientable compact surfaces, only the torus $\mathbb{R}^2/\mathbb{Z}^2$ admits expanding transformations!

ex: Give examples of expanding transformations of \mathbb{R} , of \mathbb{R}/\mathbb{Z} and of $\mathbb{R}^2/\mathbb{Z}^2$.

ex: Can an expanding map of a compact space be an homeomorphism? The answer is yes if the space is finite, and an example is not so difficult. On the other side, one can show (but it is not easy!) that an infinite compact space does not admits expanding homeomorphisms.

Decimal expansion. The most famous expanding map is of course “multiplication by 10”, the circle map $E_{10} : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ defined by

$$E_{10}(x + \mathbb{Z}) = 10 \cdot x + \mathbb{Z}$$

If $x = 0.x_1x_2x_3\dots$, with $x_n \in \{0, 1, 2, \dots, 9\}$ is the representation of $x \in [0, 1)$ in base 10, then

$$E_{10}(0.x_1x_2x_3\dots + \mathbb{Z}) = 0.x_2x_3x_4\dots + \mathbb{Z}$$

Periodic and pre-periodic points, which correspond to rationals, are dense.

If the circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ is equipped with the standard metric, it is clear that if $0 < d(x, x') < 1/20$ then $d(E_{10}(x), E_{10}(x')) = 10 \cdot d(x, x')$. Therefore, E_{10} is expanding.

Sensitive dependence on initial conditions can also easily recognized. Indeed, if $0 < d(x, x') < 1/2 \cdot 10^{-n}$, then $d(E_{10}^n(x), E_{10}^n(x')) = 10^n \cdot d(x, x')$. Therefore, for all $\varepsilon > 0$ and all $x \in \mathbb{R}/\mathbb{Z}$, there exist another point $x' \in \mathbb{R}/\mathbb{Z}$ and a time $n \geq 0$ such that

$$d(x, x') < \varepsilon \quad \text{e} \quad d(E_{10}^n(x), E_{10}^n(x')) > 1/4.$$

It is also clear that for any not-empty interval $I \subset \mathbb{R}/\mathbb{Z}$, there exists a time $n \geq 0$ such that $E_{10}^k(I) = \mathbb{R}/\mathbb{Z}$ for all times $k \geq n$ (it is sufficient to observe that I contains some interval $J = [k/10^2, (k+1)/10^n]$ for n sufficiently large, and that $f^n(J) = \mathbb{R}/\mathbb{Z}$). Thus, E_{10} is topologically mixing.

Let $\alpha = \alpha_1\alpha_2\dots\alpha_n$ be a finite word in the letters of the alphabet $\{0, 1, 2, \dots, 9\}$. There exists a residual set of points $x \in \mathbb{R}/\mathbb{Z} \approx [0, 1)$ such that their base 10 representation contains the word α infinitely often (in the sense that, if $x = 0.x_1x_2x_3\dots$, there exist an infinity of times $k \geq 0$ such that $x_{k+1}x_{k+2}\dots x_{k+n} = \alpha_1\alpha_2\dots\alpha_n$). Moreover, since finite words are countable, there exists a residual set of points $x \in \mathbb{R}/\mathbb{Z} \approx [0, 1)$ such that their base 10 representation contains all finite words in the alphabet $\{0, 1, 2, \dots, 9\}$ infinitely often. This means that a “generic” infinite book contains all the possible finite books infinitely often! More is true, as showed by Émile Borel (see the paragraph on normal numbers below).

ex: Give example of points in the residual sets described above.

ex: The *Champernowne constant* is

$$c = 0.123456789101112131415161718192021222324 \dots$$

Is the orbit of $c + \mathbb{Z}$ under E_{10} dense in the circle?

Linear expanding maps of the circle. There is, of course, nothing special with the number 10 used above, the number of fingers in our hands. A *standard expanding map* of the interval is a map $E_N : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$, defined by

$$x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$$

where $N \in \mathbb{Z}$ an integer such that $|N| > 1$. It is expanding, topologically mixing and has sensitive dependence on initial conditions, and have dense and countable set of periodic points. Proofs are just a rewriting of the above proofs that we gave for $N = 10$.

The map f is a factor of the Bernoulli shift over an alphabet made of N letters, and the set where the semi-conjugation fails to be one-to-one is small (it is made of rationals).

Together with periodic orbits and dense infinite orbits, such maps also admit more complicated orbit closures. For example, the expanding map with $N = 3$ clearly preserves the Cantor set K , thought as a subset of the circle (i.e. with the points 0 and 1 identified), i.e. $E_N(K) \subset K$, and the restriction $f|_K : K \rightarrow K$ is clearly topologically mixing. Thus, there exist orbits of $E_N : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ which are dense in K .

Non-linear expanding maps of the circle. We now consider a generic, not necessarily linear, expanding map $g : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ of class \mathcal{C}^1 , i.e. such that any of its lifts $G : \mathbb{R} \rightarrow \mathbb{R}$ is continuously differentiable. Being G' periodic with period 1, there exists $\lambda > 1$ such that $|G'(x)| > \lambda$ for all $x \in \mathbb{R}$, and G' does not change sign. In particular, the degree of g has absolute value strictly bigger than one, since

$$|\deg(g)| = |G(1) - G(0)| = \left| \int_0^1 G'(x) dx \right| = \int_0^1 |G'(x)| dx > \int_0^1 \lambda dx > 1.$$

Theorem 9.2. Any expanding map $g : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ of class \mathcal{C}^1 and degree $\deg(f) = N$ is topologically conjugated to the standard expanding map $E_N : x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$.

Proof. For simplicity, we assume that G is increasing, i.e. that $N > 1$. The idea is to first define a conjugation between the pre-images of a fixed point, and then extend it to the whole circle using the facts that such pre-images are dense.

Let $x_k^i = i/\lambda^k$, with $i = 0, 1, \dots, \lambda^k - 1$. Then $E_N(x_k^i) = x_{k-1}^{i'}$, where i' is the unique integer between 0 and $N^{k-1} - 1$ such that $i = i' \bmod N^{k-1}$. Let p be the fixed point of G , a lift of g . Since G is strictly increasing and $G(p+1) = p+N$, there exist $p = y_1^0 < y_1^1 < \dots < y_1^{N-1} < p+1$ such that $g(y_1^i) = p+i$. Inductively (in k) we define the points y_k^i , with $i = 0, 1, \dots, N^k - 1$ such that

$$y_{k-1}^i = y_k^{Ni} < y_k^{Ni+1} < \dots < y_k^{Ni+N-1} < y_k^{Ni+N} = y_{k-1}^{i+1}$$

and $G(y_k^i) = y_{k-1}^{i'}$, where i' is the unique integer between 0 and $N^{k-1} - 1$ such that $i = i' \bmod N^{k-1}$. For any interval $I_k^i = \pi([y_k^i, y_k^{i+1}])$ we have $g^k(I_k^i) = \mathbb{R}/\mathbb{Z}$ (remember that $\pi : \mathbb{R} \rightarrow \mathbb{R}/\mathbb{Z}$ is the projection of the line onto the circle). Since g is expanding, i.e. there exists $\lambda > 1$ such that $|G'(x)| > \lambda$ for all x , any of these intervals has length $|I_k^i| < \lambda^{-k}$, and therefore the family of points $\{y_k^i\}_{k \in \mathbb{N}, i=0,1,\dots,N^k-1}$ is dense in $[p, p+1]$. The function

$$H : \{y_k^i\}_{k \in \mathbb{N}, i=0,1,\dots,N^k-1} \rightarrow \{x_k^i\}_{k \in \mathbb{N}, i=0,1,\dots,N^k-1}$$

defined by $H(y_k^i) = x_k^i$ is strictly monotone. The density of the points $\{y_k^i\}$ and $\{x_k^i\}$ allows to extend H as a homeomorphism $H : [p, p+1] \rightarrow [0, 1]$, which in turn defines a homeomorphism $h : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$. Finally, one easily see that $E_N \circ h = h \circ g$. \square

In particular, given an expanding map of the circle $g : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ of class \mathcal{C}^1 , all maps sufficiently near to g in the \mathcal{C}^1 topology is topologically conjugated to E_N . This follows from the fact that expansiveness is an open condition, and that the degree is locally constant. Therefore,

Theorem 9.3. *Continuously differentiable expanding maps of the circle are \mathcal{C}^1 -structurally stable.*

9.4 Symbolic dynamics and coding

Bernoulli shift. The abstract archetypal mixing map is the *(one-sided) Bernoulli shift*

$$\sigma : \Sigma^+ \rightarrow \Sigma^+$$

where $\Sigma^+ = \mathcal{A}^{\mathbb{N}}$ is the space of infinite one-sided words $x = x_1x_2x_3\ldots$ in the finite alphabet $\mathcal{A} \approx \{1, 2, \dots, N\}$ made of $N \geq 2$ letters, and

$$\sigma(x_1x_2x_3\ldots) = x_2x_3x_4\ldots$$

We recall that, as discussed in section 5.2, the product topology in Σ^+ is metrizable. It may be induced by the ultrametric d_∞ , and open and closed balls for this metric are centered cylinders C_α , made of infinite words which start with a fixed finite word $\alpha = \alpha_1\alpha_2\ldots\alpha_k \in \mathcal{A}^k$. The larger is the (length of the) word α , the smaller is the (radius of the) ball C_α .

Any open not-empty set $U \subset \Sigma^+$ contains some centered cylinder $C_\alpha \subset U$, and, if $|\alpha|$ denotes the length of the finite word α , it is clear that $\sigma^n(C_w) = \Sigma^+$ for all times $n \geq |\alpha|$. A fortiori, $\sigma^n(U)$ intersect any other not-empty open set for such large times n . Thus, σ is topologically mixing. Being mixing, it is +transitive, and therefore a generic point has a dense orbit.

Indeed, in this case it is quite simple to exhibit points with dense orbits. Since the set of finite words in the alphabet is countable, we may enumerate finite words as $\alpha_1, \alpha_2, \alpha_3, \dots$, and then observe that the trajectory of the point $x = \alpha_1\alpha_2\alpha_3\ldots$ passes through all centered cylinders, hence through all not-empty open sets.

Less obvious is to give example of points x such that $\omega_\sigma(x) = \Sigma^+$, which is also a generic property. An example is

$$x = \alpha_1\alpha_1\alpha_2\alpha_1\alpha_2\alpha_3\alpha_1\alpha_2\alpha_3\alpha_4\alpha_1\alpha_2\alpha_3\alpha_4\alpha_5\ldots$$

whose trajectory visits all centered cylinders infinitely often

As we already saw, the Bernoulli shift also have dense periodic points, since any cylinder C_α contains the periodic point $\alpha\alpha\alpha\ldots$. In particular, the Bernoulli shift is chaotic. Also, the fixed points of σ^n have cardinality $P_n(\sigma) = \text{card}(\text{Per}_n(\sigma)) = N^n$.

As for expanding maps, there are points whose orbit is dense in a proper subset of Σ^+ . For example, the restriction of the shift on $(\mathcal{A} \setminus \{k\})^{\mathbb{N}} \subset \Sigma^+$, formed by infinite words which do no use the letter $k \in \mathcal{A}$ (or any other letter), is topologically mixing (we may just repeat the discussion above). Therefore, a generic point $x \in (\mathcal{A} \setminus \{k\})^{\mathbb{N}}$ has an orbit which is dense inside $(\mathcal{A} \setminus \{k\})^{\mathbb{N}}$. More examples are given by “subshifts”, defined below.

ex: Give examples of points $x \in \Sigma^+$ such that $\omega_\sigma(x) = \Sigma^+$.

ex: Give example of points $x \in \Sigma^+$ tais such that the closure of the orbit $O_\sigma^+(x)$ is a proper and infinite subset of Σ^+ .

ex: Give example of points $x \in \Sigma^+$ which are not recurrent.

Full shift. Also interesting is the invertible version. Let $\Sigma = \mathcal{A}^{\mathbb{Z}}$ be the space of bi-infinite words $x = \ldots x_{-2}x_{-1}x_0x_1x_2\ldots$ in the letters of the finite alphabet $\mathcal{A} = \{1, 2, \dots, \Omega\}$, equipped with the product topology. The (full) shift $\sigma : \Sigma \rightarrow \Sigma$, defined by $(\sigma(x))_k = x_{k+1}$, is a homeomorphism.

ex: Show that periodic points of the full shift are dense.

ex: Show that the full shift is topologically mixing.

Subshifts. The restriction of the shift σ to a closed invariant subset X of Σ or Σ^+ is called *subshift*, or also *symbolic dynamical system*.

Given any family \mathcal{F} of finite words, finite or not, one can define a subset $\Sigma_{\mathcal{F}} \subset \Sigma$ as the set of those infinite words which do not “contain” any of the finite words $\phi \in \mathcal{F}$ (an infinite word $x = x_1x_2x_3\ldots$ contains the finite word $\phi = \phi_1\phi_2\ldots\phi_k$ if there exists a time $n \geq 0$ such that $x_{n+1}x_{n+2}\ldots x_{n+k} = \phi_1\phi_2\ldots\phi_k$). It is clear that $\Sigma_{\mathcal{F}}$ is σ -invariant. Also, since cylinders are clopen balls, such $\Sigma_{\mathcal{F}}$ is also closed, being an intersection of closed sets. Indeed, it is easy to see that any closed invariant subset of Σ is of this type: it is defined by the family of “forbidden” words. Similarly, one defines $\Sigma_{\mathcal{F}}^+$.

If the set \mathcal{F} of forbidden words is finite, then the restriction of σ to $\Sigma_{\mathcal{F}}$ or $\Sigma_{\mathcal{F}}^+$ is called (full or one-sided) *subshift of finite type*.

Topological Markov chains. The simplest way to produce invariant subsets uses transition matrices, an idea which comes from the theory of Markov chains in probability. Let $A = (a_{ij})$ be a “transition matrix”, i.e. a $N \times N$ matrix with entries 0 or 1. Let

$$\Sigma_A^+ := \{x = x_1x_2x_3\ldots \in \Sigma^+ \text{ such that } a_{x_nx_{n+1}} = 1 \ \forall n \geq 0\}.$$

It is clear that $\sigma(\Sigma_A^+) \subset \Sigma_A^+$. The restriction

$$\sigma_A := \sigma|_{\Sigma_A^+} : \Sigma_A^+ \rightarrow \Sigma_A^+$$

is called *topological Markov chain*. The idea is that the letters of the alphabet represent the possible *states* of a system, and transition from the state $x_n = i$ at time n to the state $x_{n+1} = j$ at time $n+1$ is possible iff $a_{ij} = 1$. Thus, a topological Markov chain is a subshift of finite type, where the set of forbidden words is made of words of length two, namely $\mathcal{F} = \{ij \text{ s.t. } a_{ij} = 0\}$. In order to avoid trivialities, as for example empty sets, it is a good idea to consider only transition matrices with at least one 1 in each row, so that there is always some letter to write after any given letter.

A finite word (also called *block* by probabilists) $\alpha = \alpha_1\alpha_2\ldots\alpha_n$ is *admissible* if $a_{\alpha_k\alpha_{k+1}} = 1$ for all $k = 1, 2, \dots, n-1$, i.e. if it is a piece of an infinite word of Σ_A^+ . The set of admissible finite words is also called *language*, with an obvious analogy.

The relative topology on Σ_A^+ is generated by the intersections of centered cylinders $C_\alpha \subset \Sigma^+$ with Σ_A^+ , which are empty if α is not an admissible word. With abuse of language we will still denote by C_α such non-empty intersections, and call them “admissible cylinders”.

Counting words. Let $\sigma_A : \Sigma_A^+ \rightarrow \Sigma_A^+$ be a topological Markov chains defined by the transition matrix A . Here we count admissible words and therefore periodic points.

It is useful to introduce a *Markov graph* \mathcal{G}_A , whose vertices are the letters/states of the alphabet $\mathcal{A} \approx \{1, 2, \dots, N\}$, with oriented edges from i to j whenever $a_{ij} = 1$. Admissible words are therefore walks/paths, i.e. finite sequences of consecutive letters joined by the edges of the graph.

Let $W_n(ij)$ be the cardinality of admissible words of length $n+1$ which start with the letter i and end with the letter j (we omit any reference to the transition matrix to simplify the notation). This is the number of different walks/paths of length n in \mathcal{G}_A joining i to j . The key observation is the following.

Theorem 9.4. *The number $W_n(ij)$ is equal to the ij -entry of the n -th power of the transition matrix, i.e.*

$$W_n(ij) = (A^n)_{ij}.$$

Proof. This is obvious when $n = 1$, since $W_1(ij) = a_{ij}$, and follows easily by induction in the general case. Indeed, to get an admissible word of length $n+1$ starting from i and ending with

j , we must attach, to an admissible word of length n starting with i and ending with some k , one more letter j , and this is possible provided $a_{kj} = 1$. By the inductive hypothesis, this number is

$$W_n(ij) = \sum_k W_{n-1}(ik) W_1(kj) = \sum_k (A^{n-1})_{ik} a_{kj} = (A^n)_{ij}.$$

□

In particular, when $i = j$ we are counting the fixed points of σ_A^n . Hence, summing over all i 's, we get

Theorem 9.5. *The cardinality of n -periodic points of a topological Markov chain defined by the transition matrix A is*

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

Asymptotic growth of walks. Also interesting is estimate the asymptotic growth of the number of walks in the graph. Let W_n be the cardinality of admissible words of length $n + 1$, i.e. of walks of length n in the directed graph \mathcal{G}_A . As follows from 9.4, this number is

$$W_n = \sum_{ij} W_n(ij) = \sum_{ij} (A^n)_{ij} = \|A^n\|_1$$

(since powers of the transition matrix also have non-negative entries, thus this sums coincides with the 1-norm of the matrix A^n). Also, since the 1-norm in the algebra of square matrices is sub-multiplicative, we see that

$$W_{n+m} \leq W_n \cdot W_m.$$

Thus, these numbers grow at most exponentially, since their logarithms are sub-additive. This implies that we may define an exponential rate of growth (of the Markov chain defined by A), according to

$$w(\sigma_A) := \lim_{n \rightarrow \infty} \frac{1}{n} \log W_n$$

This number may be estimated using the Perron-Frobenius theorem, which implies that W_n grows like $\sim \rho^n$, if ρ is the eigenvalue of A with largest absolute value. Alternatively, one can use Gelfand's formula and deduce that w is the spectral radius of A , i.e. the maximal modulus of its eigenvalues. It turns out that this number is also the “topological entropy” of the topological Markov chain, to be defined later.

Irreducible and transitive Markov chains. The topological Markov chain $\sigma_A : \Sigma_A^+ \rightarrow \Sigma_A^+$ is *irreducible* if for any two states i and j there exist a time $n \geq 1$ (depending on i and j) such that the ij entry of the n -th power of A is not zero, i.e. $(A^n)_{ij} \neq 0$. This means that there exist admissible words starting from any i and ending to any j . Equivalently, this means that the graph \mathcal{G}_A is “strongly connected”, i.e. it admits paths joining any two vertices i and j (but the length of these paths depend on the vertices).

Theorem 9.6. *An irreducible topological Markov chain is topologically transitive.*

Proof. Let C_α and C_β be two admissible cylinders, defined by the admissible words $\alpha = \alpha_1 \alpha_2 \dots \alpha_i$ and $\beta = \beta_1 \beta_2 \dots \beta_j$. Since A is irreducible, there exist a time n such that the α_i - β_1 entry of the matrix A^n is non-zero. This means that there exists an admissible word γ , of length $n - 1$ (possibly empty if $n = 1$), such that $\alpha_i \gamma \beta_1$ is an admissible word. But then also $\alpha \gamma \beta$ is admissible. Since clearly $C_{\alpha \gamma \beta} \subset C_\alpha$ and $\sigma^{n+i-1}(C_{\alpha \gamma \beta}) = C_\beta$, we see that $\sigma^{n+i-1}(C_\alpha) \cap C_\beta$ contains C_β , hence it is not empty. □

The topological Markov chain $\sigma_A : \Sigma_A^+ \rightarrow \Sigma_A^+$ is *transitive* (or *aperiodic*) if there exists a time $n \geq 1$ such that all the entries of A^n are strictly positive. Thus, any two states i and j can be joined by an admissible path of length n . It is clear that this implies that each row and each column of A has at least a non-zero entry, for otherwise $A^n = AA^{n-1} = A^{n-1}A$ would have a zero entry (if you think about admissible paths this is obvious). This implies, by induction, that all the entries of A^m are strictly positive whenever $m \geq n$. Therefore, any two states of a transitive topological Markov chain are the initial and final letter of admissible words of any length $\geq n+1$. Equivalently, any two vertices i and j can be joined by a path in \mathcal{G}_A of any length $\geq n$.

Clearly, a transitive Markov chain is irreducible, but the converse is false.

Theorem 9.7. *A transitive topological Markov chain is topologically mixing and has dense periodic points.*

Proof. The same proof as above, with a uniform value of n which works for all possible pairs $\alpha_n\beta_1$, shows that a transitive Markov chain is mixing. Moreover, any admissible cylinder C_α contains an admissible cylinder $C_{\alpha\beta}$ with some admissible word $\alpha\beta$ of length n as above, and such cylinder contains the admissible periodic point $\alpha\beta\alpha\beta\alpha\beta\dots$ \square

ex: The Markov chain defined by the matrix E with all entries equal to one is the Bernoulli shift. Compute the quantities $W_n(ij)$, W_n , and $P_n(\sigma_E)$ and $w(\sigma_E)$.

ex: Consider the topological Markov chains defined by the transition matrices

$$A = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \quad B = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \quad C = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \quad D = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Compute the number P_n of n -periodic points and the number W_n of different words of length $n+1$. Are these topological Markov chains irreducible? Are they transitive?

ex: Prove (along the lines suggested above) that if all the entries of the n -th power A^n of the transition matrix A are positive, then also all the entries of A^m are positive, whenever $m \geq n$.

Prime walks theorem. Consider a topological Markov chain $\sigma_A : \Sigma_A^+ \rightarrow \Sigma_A^+$ defined by the transition matrix A . Since $P_n = P_n(\sigma_A)$ counts the number of closed walks of length n in the directed graph \mathcal{G}_A , it can be regarded as a geometric object. Cyclic permutations on the vertices of a closed walk produce other “geometrically equivalent” walks, thus define the same closed curve p of length $|p| = m$ such that $m \mid n$. We may therefore consider the cardinality

$$\Pi_n := \text{card}\{\text{closed curves } p \text{ s.t. } |p| = n\}$$

of *prime* closed walks of length n . Thus,

$$P_n = \sum_{m \mid n} m \Pi_m$$

and therefore, by Möbius inversion formula (see for example [HW59] theorem 266),

$$n \Pi_n = \sum_{m \mid n} \mu(m) P_{n/m} \quad (9.2)$$

where $\mu(n)$ is the Möbius function (which is equal to 1 if n is a square-free integer with an even number of prime factors, to -1 if n is a square-free integer with an odd number of prime factors, and zero otherwise). The Perron-Frobenius theorem 9.9 implies that a transitive transition matrix A has a positive simple eigenvalue λ_1 with is strictly larger than the modulus $|\lambda_k|$ of all

other eigenvalues $\lambda_2, \lambda_3, \dots$. In particular, the spectral radius of A is $\rho = \lambda_1$. There follows from theorem 9.5 that

$$P_n = \text{tr} A^n = \rho^n + o(\rho^n)$$

for large n . The leading term in the sum (9.2) is when $m = 1$, and others are smaller, since

$$n \Pi_n = P_n + \sum_{m|n, m \geq 2} \mu(m) P_{n/m} = \rho^n + \mathcal{O}(\rho^{n/2}).$$

Finally, the “prime walks theorem” can be stated as

Theorem 9.8. *Let Π_n be the cardinality of prime closed walks in the graph defined by a transitive transition matrix A with spectral radius ρ . Then*

$$\Pi_n \sim \frac{\rho^n}{n}$$

It is evident the analogy with the “prime number theorem”, which says that the cardinality of prime numbers smaller than x is asymptotic to $\pi(x) \sim x/\log x$. Deep generalizations have been found by Parry and Pollicott^{38 39} and use, as expected, sophisticated methods of analytic number theory, starting from zeta functions.

Coding. Symbolic dynamical systems are abstract models for dynamical systems. One of the central idea in dynamical systems is indeed to “code” an actual map $f : X \rightarrow X$ with a symbolic system.

A possible strategy is the following. We divide the phase space X into a finite number of disjoint pieces

$$X = B_1 \cup B_2 \cup \dots \cup B_N$$

The trajectory of any point $x \in X$ visits the different pieces according to some pattern, hence defines an “itinerary”, which is the infinite word $h(x) = x_0 x_1 x_2 \dots$ of letters in the alphabet $\mathcal{A} = \{1, 2, \dots, N\}$ defined according to $x_n = k$ iff $f^n(x) \in B_k$.

We define the transition matrix $A = (a_{ij})$ such that $a_{ij} = 1$ if $B_j \subset f(B_i)$ and $a_{ij} = 0$ otherwise. It is clear that possible histories of points of X belong to Σ_A^+ , so that $h(X) \subset \Sigma_A^+$. Since $f^n(f(x)) = f^{n+1}(x)$, if $h(x) = x_0 x_1 x_2 \dots$ is the itinerary of x , then the itinerary of $f(x)$ is

$$h(f(x)) = x_1 x_2 x_3 \dots = \sigma_A^+(x_0 x_1 x_2 \dots) = \sigma_A^+(h(x)).$$

That is, the coding map intertwines between f and σ_A^+ , i.e.

$$\sigma_A \circ h = h \circ f.$$

If it happens to be surjective, it defines a semi-conjugation. In general, the coding map h is neither injective nor surjective. One way to get more admissible itineraries is to include boundaries (when this makes sense) into the definition of the B_k ’s, thus allowing not-empty (but of zero measure) intersections between the pieces of the “partition”.

An alternative, is to look for a semi-conjugation in the opposite direction. If f is sufficiently chaotic, one may hope that to any itinerary $x_0 x_1 x_2 \dots \in \Sigma_A^+$ there corresponds a unique point

$$\{x\} = \bigcap_{n=0}^{\infty} f^{-n}(B_{x_n})$$

of the phase space. This would give a map $\ell : \Sigma_A^+ \rightarrow X \dots$

³⁸W. Parry, An analogue of the prime number theorem for shifts of finite type and their suspensions, *Israel Journal of Mathematics* **45** (1983), 41-52.

³⁹W. Parry and M. Pollicott, An analogue of the prime number theorem for closed orbits of axiom A flows, *Annals of Mathematics* **118** (1983), 573-591.

Binary expansion. Consider the multiplication by two map $f : [0, 1] \rightarrow [0, 1]$, defined by $f(x) = \{2x\}$. The natural partition is given by $B_0 = [0, 1/2)$ and $B_1 = [1/2, 1]$. Since $f(B_0) = [0, 1]$ and $f(B_1) = [0, 1]$, the transition matrix has all entries equal to one: all transitions are allowed. The itinerary of any point $x \in [0, 1]$ defines a sequence $x_1 x_2 x_3 \dots$ of 0's and 1's such that $f^n(x) \in B_{x_n}$. It is clear that we can recover the point as

$$x = 0.x_1 x_2 x_3 \dots = \frac{x_1}{2} + \frac{x_2}{2^2} + \frac{x_3}{2^3} + \dots$$

Thus, the itinerary is a binary representation of the point (one of the two, if the point is rational). Vice-versa, to any sequence $x_1 x_2 x_3 \dots$ of 0's and 1's we may associate the number $x = 0.x_1 x_2 x_3 \dots = \sum_{k=1}^{\infty} x_k / 2^k$.

Golden ratio shift. Consider the *golden ratio map* $f : [0, 1] \rightarrow [0, 1]$, sending x to the fractional part of γx , i.e. $f(x) := \{\gamma x\}$, where $\gamma = (1 + \sqrt{5})/2 \simeq 1.618\dots$ is the Greeks' ratio. Consider the partition of the unit interval given by $B_0 = [0, 1/\gamma]$ and $B_1 = [1/\gamma, 1]$. Since the ratio is a root of the quadratic equation $\gamma = 1 + 1/\gamma$, it follows that $f(B_0) = [0, 1]$ and $f(B_1) = B_0$. The transition matrix is therefore

$$G = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

The corresponding space of sequences Σ_G^+ is made of sequences in the letters 0 and 1 with no consecutive 1's, i.e. which do not contain the word 11. It is known as *golden ratio shift*.

One can show that the point corresponding to the itinerary $x_1 x_2 x_3 \dots$ is

$$x = x_1 \gamma^{-1} + x_2 \gamma^{-2} + x_3 \gamma^{-3} + \dots$$

(such expansions using non-integer bases are called *beta expansions*).

ex: Show that the number W_n of admissible words of length n for the golden ratio shift satisfies the Fibonacci recursion

$$W_{n+2} = W_{n+1} + W_n.$$

Conclude that it grows like $W_n \sim C \gamma^n$ for some constant C and large n .

ex: The *Backer's map* is the transformation $f : [0, 1]^2 \rightarrow [0, 1]^2$ of the unit square defined by

$$(x, y) \mapsto \begin{cases} (2x, y/2) & \text{if } 0 \leq x \leq 1/2 \\ (2x - 1, (y + 1)/2) & \text{if } 1/2 < x \leq 1 \end{cases}$$

Discuss its dynamics. Consider the full shift $\sigma : \Sigma \rightarrow \Sigma$ on $\Sigma = \{0, 1\}^{\mathbb{Z}}$. Show that the map $h : \Sigma \rightarrow [0, 1]^2$, defined by

$$\dots x_{-2} x_{-1} x_0 x_1 x_2 \dots \mapsto \left(\sum_{n=0}^{\infty} \frac{x_{-n}}{2^n}, \sum_{n=1}^{\infty} \frac{x_n}{2^n} \right)$$

is a semi-conjugation between σ and f .

9.5 Non-negative matrices and the Perron-Frobenius theorem

Non-negative and stochastic matrices. Transition matrices belong to the larger class of *non-negative matrices*, those square $n \times n$ real matrices $A = (a_{ij})$ with non-negative entries $a_{ij} \geq 0$ for all $1 \leq i, j \leq n$ (nothing to do with “positive-defined” symmetric matrices!). Non-negative matrices define the linear maps of \mathbb{R}^N that preserve the closure of positive cone \mathbb{R}_+^N , the convex cone of non-negative vectors, made of those $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^N$ such that $x_k \geq 0$ for all k .

Non-negative matrices also contain the class of *stochastic matrices*, those non-negative matrices $P = (p_{ij})$ such that each column is a probability vector, i.e. $\sum_i p_{ij} = 1$ for all j . They are precisely the non-negative matrices that preserve the unit simplex $\Delta^N \subset \mathbb{R}^N$, the convex set of probability vectors, those $p \in \mathbb{R}^N$ such that $0 \leq p_k \leq 1$ and $\sum_k p_k = 1$.

It is clear that the definitions of irreducible and transitive also make sense for non-negative matrices. Non-negative matrices are the characters of the celebrated and useful Perron-Frobenius theorem^{40 41 42}, which says, in one of its versions and between other things, that their spectral radius is a simple eigenvalue with a positive eigenvector.

Theorem 9.9 (Perron-Frobenius). *Let A be a transitive non-negative square matrix. Then there exists a simple positive eigenvalue ρ , with an associated eigenvector which is also positive (i.e. has strictly positive coordinates), and all other eigenvalues λ have strictly smaller modulus, i.e. $|\lambda| < \rho$.*

An elementary proof may be found in [KH95]. Existence of the positive/dominant eigenvector may be shortened using the Brouwer fixed point theorem 6.2, and therefore at the expense of some non-trivial algebraic geometry. Beautiful geometric proofs (of existence of a positive eigenvector) were later found by Samelson⁴³ and Birkhoff⁴⁴ (Garrett, the son of George David), using the Hilbert projective metric in the interior of a cone and allowing generalizations in infinite dimension. These ideas gave origin to the modern theory of “Perron-Frobenius”, or “transfer”, operators, thanks to David Ruelle, . . .

Since the Perron-Frobenius theorem is important also in the ergodic theory of Markov chains (as well as in many application to modern technology . . .), it is worth some paragraphs. We follow Birkhoff’s original exposition.

The hyperbolic line. Consider the cone \mathbb{R}_+^2 made of vectors $x = (x_1, x_2) \in \mathbb{R}^2$ with non-negative coordinates $x_k \geq 0$. Let \mathbb{H} be the space of rays of \mathbb{R}_+^2 , or “projective cone”, whose points are rays $x^+ = \{tx : t > 0\}$ through non-zero vectors x of \mathbb{R}_+^2 . It can be naturally identified with the unit simplex Δ^2 , the segment between the basis vectors $(0, 1)$ and $(1, 0)$, made of non-negative vectors $x = (x_1, x_2)$ such that $x_1 + x_2 = 1$. The space of rays may also be given homogeneous coordinate $\xi = x_2/x_1$, sending the vertices $(1, 0)$ and $(0, 1)$ of the simplex to 0 and ∞ , respectively. The map $x \simeq \xi \mapsto \log \xi = \log(x_2/x_1)$ sends therefore the projective cone into the extended real line $[-\infty, \infty]$. We define the “hyperbolic metric” in the interior of \mathbb{H} demanding that this map is an isometry, namely

$$\rho(\xi, \eta) := |\log \xi - \log \eta| = |\log(\xi/\eta)|$$

if $\xi = x_2/x_1$ and $\eta = y_2/y_1$. The boundary points $\xi = 0$ and $\xi = \infty$ are, by definition, at infinite distance from all other points.

Recall the the *cross-ratio* between the four ordered points $\alpha < \xi < \eta < \beta$ in the extended line is

$$R(\alpha, \xi, \eta, \beta) := \frac{(\eta - \alpha)(\beta - \xi)}{(\xi - \alpha)(\beta - \eta)}$$

It is clear that the cross-ratio is invariant under homotheties $\xi \mapsto a\xi$, with $a \neq 0$, and translations $\xi \mapsto \xi + b$, and an elementary computation shows that it is also invariant under the inversion $\xi \mapsto 1/\xi$ (provided we interpret $(\infty - \xi)/(\infty - \eta) = 1$, as convenient). There follows that the cross-ratio is invariant under all projective transformations $\xi \mapsto f(\xi) = (a\xi + b)/(c\xi + d)$, with $\alpha d - \beta\gamma \neq 0$, namely, $R(f(\alpha), f(\xi), f(\eta), f(\beta)) = R(\alpha, \xi, \eta, \beta)$.

Thus, we recognise that the hyperbolic metric in the interior of the simplex is

$$\rho(\xi, \eta) = |\log R(0, \xi, \eta, \infty)| .$$

Conversely, the line $t \mapsto (1 - t)x + ty$ cuts the convex cone \mathbb{R}_+^2 in a bounded interval between $t = \alpha$ and $t = \beta$, which intersects the rays through x and y at times $t = 0$ and $t = 1$, respectively. There exists a projective transformation f , unique up to homotheties, that sends $f(\alpha) = 0$ and $f(\beta) = \infty$. Then the hyperbolic distance between the rays through x and y is $|\log f(0)/f(1)| = |\log R(0, f(0), f(1), \infty)| = |\log R(\alpha, 0, 1, \beta)|$.

⁴⁰O. Perron, Grundlagen für eine Theorie des Jacobischen Kettenbruchalgorithmus, *Math. Ann.* **64** (1907), 11-76.

⁴¹O. Perron, Zur Theorie de Matrices, *Math. Ann.* **64** (1907), 248-263.

⁴²G. Frobenius, Über Matrizen aus positiven Elementen I, II, *S.-B. kgl. Preuss. Acad. Berlin* (1908), 471-476; (1909), 514-518.

⁴³H. Samelson, On the Perron-Frobenius theorem, *Michigan Math.* **4** (1956), 57-59.

⁴⁴G. Birkhoff, Extensions of Jentzsch’s theorem, *Trans. Amer. Math. Soc.* **85** (1957), 219-227.

A linear transformation of the plane

$$(x_1, x_2) \mapsto (ax_2 + bx_1, cx_2 + dx_1)$$

defined by a 2×2 matrix with non-negative entries $a, b, c, d \geq 0$ clearly preserves the cone \mathbb{R}_+^2 , and sends rays into rays. Therefore, it induces a projective map $f : \mathbb{H} \rightarrow \mathbb{H}$, which is $f(\xi) = (a\xi + b)/(b\xi + d)$ in homogeneous coordinates. We want to compute its Lipschitz constant w.r.t. the hyperbolic metric. It is clear that the inversion $\xi \mapsto 1/\xi$ is an isometry, but a generic projective map is not. The infinitesimal hyperbolic distance between the points $\xi + d\xi$ and ξ is $|\log((\xi + d\xi)/\xi)| = |\log(1 + d\xi/\xi)| \simeq |d\xi/\xi|$. An elementary computation shows that the Lipschitz constant of the projective map $\xi' = f(\xi)$ is bounded by

$$\sup_{\xi} \left| \frac{\xi}{\xi'} \frac{d\xi'}{d\xi} \right| = \sup_{\xi} \left| \frac{\xi}{a\xi + b} \frac{ad - bc}{c\xi + d} \right| = \frac{\sqrt{k} - 1}{\sqrt{k} + 1}$$

where $k = ad/bc$, which we may assume positive (otherwise, reverse orientation using the isometry $\xi \mapsto 1/\xi$). In particular, such projective transformations do not increase the hyperbolic distances.

We now observe that the projective map $f(\xi) = (a\xi + b)/(c\xi + d)$ sends the origin to $f(0) = b/d$ and the point at infinity to $f(\infty) = a/b$. If these quotients are finite, then the hyperbolic diameter of the image $f(\mathbb{H})$ is precisely $\rho(a/b, c/d) = \log(ad/bc) = \log k$. There follows from the above estimate

Theorem 9.10 (Birkhoff). *A projective transformation f induced by a non-negative matrix does not increase the hyperbolic distances. Moreover, if its image has finite hyperbolic diameter equal to $\text{diam}_{\rho}(f(\mathbb{H})) = D$, then f is a contraction, i.e.*

$$\rho(f(\xi), f(\eta)) \leq \lambda \rho(\xi, \eta)$$

with Lipschitz constant bounded by $\lambda = \tanh(D/4) < 1$.

Hilbert-Birkhoff projective metric. Let V be a real vector space, a finite vector space as \mathbb{R}^N or an infinite dimensional Banach space, and let K be a convex closed cone in V .

Recall the K is a *cone* if $tK = K$ for all $t > 0$, i.e. if it contains the entire rays $x^+ = \{tx : t > 0\}$ passing from all its vectors $x \in K$. Thus, the cone is convex if it contains the entire “angle” $\angle xy = \{tx + sy : t > 0, s > 0\}$ between all pairs of its vectors $x, y \in K$, made of all rays generated by the vectors in the segment $[xy]$ between x and y . Closeness is clear in finite dimension or when V is a Banach space. If it only has a linear space structure, we may ask the following: if $x, y \in K$ and also $x - t_n y \in K$ for a sequence $t_n \rightarrow t$, then also $x - ty \in K \cup \{0\}$.

We also assume that $-K \cap K \subset \{0\}$ (the cone does not contain a non-zero vector v and its opposite v), so that it is a “genuine” cone. We denote simply by K^+ the projective cone, the space of rays in $K \setminus \{0\}$.

The intersection of any line through two points x and y of the cone that do not define the same ray is a closed and proper subinterval $\ell = [\alpha, \beta]$ of the extended real line, where the rays through x and y correspond to $\alpha \leq \xi < \eta \leq \beta$, respectively. There is a projective transformation f , unique up to homotheties, that sends ℓ onto $[0, \infty]$. The (*Hilbert-Birkhoff*) *projective metric* in the projective cone K^+ is defined as

$$\rho_K(x^+, y^+) := \rho(f(\xi), f(\eta))$$

An alternative definitions is useful. Given two vectors $x, y \in K$, we set

$$m_K(x, y) := \sup\{t > 0 \text{ s.t. } y - tx \in K\} \quad \text{and} \quad M_K(x, y) := \inf\{t > 0 \text{ s.t. } tx - y \in K\}$$

where, by definition, $\sup \emptyset = 0$ and $\inf \emptyset = \infty$. One easily shows that $m_K(x, y) \leq M_K(x, y)$, and that equality holds iff x and y belong to the same ray. Also, it is clear that $m_K(x, y) = M_K(y, x)^{-1}$ and $M_K(x, y) = m_K(y, x)^{-1}$. With little more effort, one can show that $m_K(x, y) m_K(y, z) \leq m_K(x, z)$ and $M_K(x, y) M_K(y, z) \geq M_K(x, z)$ for all $x, y, z \in K$. Finally, define

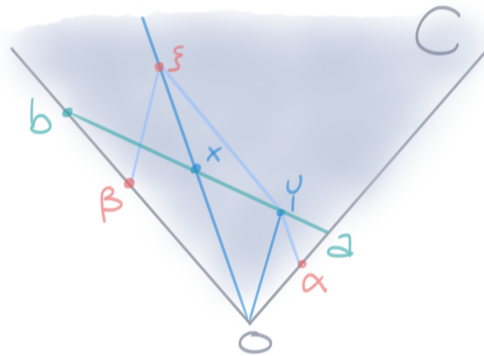
$$\rho_K(x, y) := \log M_K(x, y) - \log m_K(x, y)$$

It is clear that it only depends on the rays x^+ and y^+ . Moreover, it is non-negative, and vanishes iff x and y belong to the same ray. It follows from the previous properties that it is symmetric and satisfies the triangular inequality. It is therefore a genuine metric once restricted to rays at a finite distance from a fixed ray.

To check equivalence of the two definitions, we observe that these numbers m_K and M_K only depend on the planar section of the cone spanned by x and y , shown in the picture below. It is clear from the picture that $m_K(x, y) = |\alpha y|/|0x|$ and $M_K(x, y) = |0\xi|/|0x|$. The line through x and y intersects the cone in a closed segment ab (provided x and y are not in the same ray). Since $|\alpha y|/|0x| = |ay|/|ax|$ and $|0\xi|/|0x| = |by|/|bx|$, we see that the distance between x and y is the logarithm of the “cross-ratio” $R(b, x, y, a)$, namely

$$\rho_K(x, y) := \log \frac{|by| |ax|}{|bx| |ay|}$$

which coincides with our first definition.



It is then straightforward, from the formula using cross ratios, that if we consider a moving point $x_t = x + t(y - x)$ in the segment between x and y , with $0 \leq t \leq 1$, then $\rho_K(x, x_t) + \rho_K(x_t, y) = \rho_K(x, y)$. Therefore, segments are geodesics for the Hilbert-Birkhoff metric.

What is not clear, and indeed not true in general, is that the projective cone, equipped with this metric, is complete. This should be checked case by case.

ex: Prove that $m_K(x, y) = M_K(y, x)^{-1}$ and $M_K(x, y) = m_K(y, x)^{-1}$ (hint: if $y - tx \in K$ then $t(y/t - x) \in K$ and by convexity also $y/t - x \in K \dots$). Deduce that ρ_K is symmetric.

ex: Prove that $m_K(x, y)m_K(y, z) \leq m_K(x, z)$ and $M_K(x, y)M_K(y, z) \geq M_K(x, z)$ (hint: if $y - tx \in K$ and $z - sy \in K$ then by convexity also $sy - stx \in K$ and therefore also $z - stx \in K \dots$). Deduce that ρ_K satisfies the triangular inequality.

Linear maps and projective contraction theorem. Linear maps between real vector spaces send rays into rays. If $A : V \rightarrow W$ is a linear map sending the convex closed cone $K \subset V$ into the convex closed cone $H \subset W$, then we may consider the restriction $\varphi_A : K^+ \rightarrow H^+$ of the projective map, defined by $\varphi_A(x^+) := (Ax)^+$. It is straightforward that $m_K(x, y) \leq m_H(Ax, Ay)$ and $M_K(x, y) \geq M_H(Ax, Ay)$. Thus, the projective map does not increase the projective metric, i.e.

$$\rho_H(\varphi_A(x^+), \varphi_A(y^+)) \leq \rho_K(x^+, y^+).$$

In particular, the projective map induced by an endomorphism $A : V \rightarrow V$ which preserves the convex cone $K \subset V$ does not increase the projective metric on K . Moreover, it follows from 9.10 and our first definition of the projective metric that

Theorem 9.11 (Birkhoff-Schwartz-Pick lemma). *Let $K \subset V$ be a closed convex cone and let $A : V \rightarrow V$ be a linear map which sends the cone into itself. Then the projective map $\varphi_A : K^+ \rightarrow K^+$ does not increase the projective distances. If, moreover, the image $A(K)$ has finite projective diameter $\leq D$, then the projective map φ_A is a strict contraction, namely,*

$$\rho_K(\varphi_A(x^+), \varphi_A(y^+)) \leq \lambda \rho_K(x^+, y^+).$$

with Lipschitz constant $\lambda = \tanh(D/4)$.

If it happens that the projective cone, equipped with the Birkhoff metric, is also complete, then we may apply the contraction principle to find iteratively fixed rays for linear maps ...

e.g. Hyperbolic disk. The motivating example is the convex cone $C = \{z > x^2 + y^2\} \subset \mathbb{R}^3$, whose projectivization may be identified with the unit disk $\mathbb{D} = \{x + iy \in \mathbb{C} : x^2 + y^2 < 1\}$ of the complex plane. The Hilbert-Birkhoff metric then coincides with the hyperbolic metric of the Klein model of the hyperbolic plane, defined using cross-ratios, and such that geodesics are Euclidean segments. The above theorem 9.10 is therefore a generalization of the Schwartz-Pick lemma from complex analysis (see [Ah78]).

ex: The segment $(-1, 1)$ is the natural projectivization of the cone $C = \{y > x^2\} \subset \mathbb{R}^2$. Compute its Hilbert-Birkhoff metric.

Hyperbolic metric in the unit simplex. The example which is relevant for our proof of the Perron-Frobenius theorem is the convex cone \mathbb{R}_+^N of vectors $x = (x_1, x_2, \dots, x_N) \in \mathbb{R}^N$ with non-negative coordinates $x_k \geq 0$. Its projectivization may be naturally identified with the unit simplex Δ^N , the convex space of probability vectors, those $p = (p_1, p_2, \dots, p_N) \in \mathbb{R}_+^N$ such that $\sum_k p_k = 1$. Observe that a stochastic matrix preserves, by definition, the unit simplex, so that it may be identified with its projectivization.

The Hilbert-Birkhoff projective metric in Δ^N reads

$$\rho_\Delta(p, q) = \log M(p, q) - \log m(p, q)$$

where $m(p, q) := \min_k (p_k/q_k)$ and $M(p, q) := \max_k (p_k/q_k)$. It is clearly finite between points in the interior of the simplex, and it is clear that the interior of the simplex is complete for this metric. The boundary of the simplex is made of points at infinite distance from points at the interior (like the ideal boundary of the hyperbolic plane).

Proof of the Perron-Frobenius theorem. The hard part of the theorem, existence of the dominant eigenvalue, will be stated separately, since it is all that matter in many applications (as Markov chains, economy, ...).

A non-negative $N \times N$ matrix A clearly preserves the cone \mathbb{R}_+^N and sends rays to rays, hence induces a projective map $\varphi_A : \Delta^N \rightarrow \Delta^N$. The crucial observation is the following. The image of the unit simplex is also a simplex, the convex hull of the images (not necessarily distinct) of the vertices of the unit simplex, which are the base unit vectors e_k . Those images are rays through the columns of A , say $A_k = Ae_k$. If all the entries of A are positive, these points are in the interior of the unit simplex. It is clear that the diameter of the image $\varphi_A(\Delta^N)$ is attained as the distance between two such points. Thus,

$$\text{diam}(\varphi_A(\Delta^N)) = \max_{i,j} \rho_\Delta(A_i, A_j)$$

If all the entries, hence all the coordinates of the columns, of A are strictly positive, this diameter is clearly bounded. Birkhoff theorem 9.11 and the contraction principle therefore imply

Theorem 9.12. *Let A be a non-negative $N \times N$ matrix such that some power A^n has strictly positive entries. Then A fixes one and only one ray v^+ inside \mathbb{R}_+^N , and the generating vector v , which is strictly positive, is a eigenvector with eigenvalue $\rho > 0$. If A is a stochastic matrix then $\rho = 1$.*

Indeed, the unique fixed ray in \mathbb{R}_+^N may be found as the intersection

$$v^+ = \bigcap_{k \geq 1} A^k(\mathbb{R}_+^N)$$

which is contained in $\cap_{k \geq 1} A^{kn}(\mathbb{R}_+^N)$, which in turn is the unique fixed point of the contraction induced by A^n . This fact is important in applications, since it provides a constructive method to approximate the positive eigenvector, with exponential velocity.

To finish the proof of the Perron-Frobenius theorem 9.9 we must show that v is the only eigenvector with eigenvalue ρ , and that any other eigenvalue λ has strictly smaller absolute value $|\lambda| < \rho$. This is elementary, although laborious.

Proof. (end of the proof of theorem 9.9) We know, by theorem 9.12, that v is the only non-negative eigenvector of A , and that if w is any non-negative vector, then $A^n w$ is contained in a region of finite hyperbolic diameter inside the positive octant for all sufficiently large times n . Any other eigenvector w is necessarily outside the cone \mathbb{R}_+^N and not parallel to v .

If the eigenvalue λ is real, so that $Aw = \lambda w$, then also w is real. It is clear $v + tw \in \mathbb{R}_+^N$ if $t > 0$ is sufficiently small. But then, if $|\lambda| > \rho$, the images

$$A^{2n}(v + tw) = \rho^{2n}v + t|\lambda|^{2n}w = |\lambda|^{2n} (tw + (\rho/|\lambda|)^{2n} v) \simeq |\lambda|^{2n}tw$$

would be outside \mathbb{R}_+^N for all sufficiently large times n . Also, if $|\lambda| = \rho$, we may choose t such that $v + tw$ lies at the boundary of \mathbb{R}_+^N . Its images $A^{2n}(v + tw) = \rho^{2n}(v + tw)$ would stay at the boundary for all times n . Both phenomena cannot happen by what said at the beginning.

If the eigenvalue is not real, then it is equal to $\lambda = re^{i\theta}$ for some real $r \geq 0$ and some angle θ not multiple of π . Then (by elementary linear algebra) we may find a plane, generated by two real vectors e and f , where A acts as a rotation by an angle θ followed by a homothety of factor λ , i.e. $Ae = r(e \cos \theta - f \sin \theta)$ and $Af = r(e \sin \theta + f \cos \theta)$. But a rotation is recurrent (actually periodic if θ is rational). It is then easy to adapt the previous argument, hence find a positive or non-negative vector w in the 3-dimensional space spanned by e , f and v , such that $A^n w$ does not fall inside a region of finite hyperbolic diameter inside the positive octant if $r \geq \rho$.

A similar argument also shows that the eigenvalue ρ is simple, i.e. has algebraic multiplicity equal to one. Indeed, suppose we had a vector w such that $Aw = \rho(w + v)$. Then $v - tw$ would be positive for sufficiently small $t > 0$, but $A^n(v - tw) = \rho^n v - \rho^n t(w + nv) = \rho^n((1 - tn)v - tw)$ is certainly outside the cone for large n . \square

ex: Prove that if all the entries of a non-negative matrix A are strictly positive, say bounded below by $a_{ij} \geq \varepsilon > 0$, then the hyperbolic diameter of the image $\varphi_A(\Delta)$ is bounded by some constant $C < \infty$ (and estimate C).

ex: Estimate the hyperbolic radius of $P(\Delta)$ when P is a stochastic matrix with strictly positive entries bounded below by $p_{ij} \geq \delta > 0$.

Pagerank algorithm. A “network” made of N “nodes” and certain “links” between them may be described by a directed graph, hence by an $N \times N$ adjacency matrix $A = (a_{ij})$ with entries equal to $a_{ij} = 1$ if there is a link from the node j pointing to the node i , and $a_{ij} = 0$ otherwise. If we divide each column by the number of links from the node that it represents, we get a stochastic matrix $P = (p_{ij})$, with $p_{ij} = a_{ij} / (\sum_i a_{ij})$. Now, this matrix need not be transitive, and indeed in real world situations is a large matrix with very few non-zero entries (is a so called “sparse matrix”). Nevertheless, and this was the brilliant idea of Sergey Brin and Larry Page⁴⁵, we may add a (small) “dumping factor” $\delta > 0$ (a kind of diffusion coefficient, the probability to jump to another node chosen between all the nodes with uniform probability), and deform it to

$$G = (1 - \delta)A + \delta E$$

where E denotes the uniform stochastic matrix with all entries equal to $1/N$. This is the famous “Google matrix”. Since it is a strictly positive stochastic matrix, its iterations converge rapidly to the unique invariant probability vector p , in the interior of the positive cone, satisfying $Gp = p$. The weights p_k give good measures of the relative “popularities” of the different nodes ...

⁴⁵S. Brin and L. Page, The anatomy of a large scale hypertextual web search engine, *Computer Networks and ISDN Systems* **30** (1998), 107-117.

9.6 Cantor sets

Orbit closures of sufficiently chaotic maps have often complicated structures. If disconnected, they are typically Cantor sets, i.e. perfect and totally disconnected compact sets.

Middle-third Cantor set. The archetype is the *middle-third Cantor set*⁴⁶

$$K := \left\{ \sum_{n=1}^{\infty} \frac{x_n}{3^n} \text{ with } x_n \in \{0, 2\} \right\} \subset [0, 1],$$

the set of those numbers in the unit interval such that their base 3 representation does not use the letter “1”.

Another popular definition is $K = [0, 1] \setminus \bigcup_{k=1}^{\infty} I_k$, where the open intervals I_k are defined inductively as follows: $I_1 = (1/3, 2/3)$ is the central middle-third of the unit interval, $I_2 = (1/9, 2/9)$ and $I_3 = (7/9, 8/9)$ are the central middle-third intervals of the two components of $[0, 1] \setminus I_1$, and so on.

One more definition is $K = \bigcap_{k \geq 0} K_n$, where

$$K_n = \left\{ \sum_{k=1}^{\infty} \frac{x_k}{3^k} \text{ with } x_k \in \{0, 2\} \text{ if } k \leq n \text{ and } x_k \in \{0, 1, 2\} \text{ if } k > n \right\}$$

denotes the compact set of those numbers in the unit interval such that their base 3 representation does not use the letter “1” at the first n places. Observe that the K_n ’s form a decreasing family, i.e. $\dots \subset K_{n+1} \subset K_n \subset \dots \subset K_0 = [0, 1]$, and that each K_n is a disjoint union of 2^n closed intervals of length 3^{-n} .

In particular, K is compact and not empty being a countable intersection of a decreasing family of compact sets.

K does not contain isolated points, and therefore $K' = K$, i.e. it is “perfect”. Indeed, if $x = 0.x_1x_2x_3\dots$ is the base 3 representation of $x \in K$, we may change just the n -th digit (from 0 to 2 or vice-versa), for $n = 1, 2, 3, \dots$, and construct a sequence $x^{(n)}$ of distinct points of K converging to x .

K is “totally disconnected”, i.e. the connected component of each $x \in K$ is $\{x\}$ itself. Indeed, any two distinct point are at a distance larger than 3^{-n} for some sufficiently large n , and therefore cannot be contained in the same connected component of K_n . Thus, K is a “nowhere dense” subset of the interval, a closed subset with empty interior.

The strange properties of the Cantor sets become less mysterious if one observes that it is homeomorphic to the topological product $\{0, 2\}^{\mathbb{N}}$, the space of the Bernoulli shift over an alphabet of two letters. The homeomorphism $\varphi : \{0, 2\}^{\mathbb{N}} \rightarrow K$ is simply

$$x_1x_2\dots x_n\dots \mapsto \sum_{n=1}^{\infty} \frac{x_n}{3^n}$$

The function $\phi : \{0, 2\}^{\mathbb{N}} \rightarrow \{0, 2\}^{\mathbb{N}} \times \{0, 2\}^{\mathbb{N}}$, defined by

$$x_1x_2x_3x_4\dots \mapsto (x_1x_3\dots, x_2x_4\dots)$$

induces a homeomorphism of K onto $K \times K$. By induction, we see that K is homeomorphic to any finite power K^n . Indeed, one can prove that K is also homeomorphic to the countable Cartesian product $K^{\mathbb{N}}$ (provided one understands the product topology on this space).

Observe that $\{0, 2\}^{\mathbb{N}}$ is trivially homeomorphic to $\{0, 1\}^{\mathbb{N}}$, for example by the map sending $x_k \mapsto y_k = x_k/2$. But the binary representation defines a continuous map of $\{0, 1\}^{\mathbb{N}}$ onto the unit interval $[0, 1]$, given explicitly by $y_1y_2y_3\dots \mapsto \sum_{k=1}^{\infty} y_k/2^k$. Thus, there exists a continuous map $\psi : K \rightarrow [0, 1]$ from the Cantor set K onto the unit interval $[0, 1]$, given explicitly by

$$\sum_{n=1}^{\infty} \frac{x_n}{3^n} \mapsto \sum_{n=1}^{\infty} \frac{x_n/2}{2^n}. \quad (9.3)$$

⁴⁶G. Cantor, Über unendliche, lineare Punktmannigfaltigkeiten V, *Mathematische Annalen* **21** (1883), 545-591.

Since the Cantor set is a subset of the unit interval, by the Schröder-Bernstein theorem, K has the cardinality of the interval.

Another much appreciated property of the Cantor set is its “self-similarity”, a property which makes of K the prototype of a “fractal set”. It is clear, indeed, that any of the closed intervals which form K_n contains an affine copy of K itself (we must only make an homothety of ratio 3^{-n} and an appropriate translation).

Finally, the “length” (i.e. the Lebesgue measure) of K is

$$|K| = \lim_{n \rightarrow \infty} |K_n| = \lim_{n \rightarrow \infty} 2^n \cdot 3^{-n} = 0.$$

The Cantor set is very “small”, while containing the same cardinality of points as the whole interval!

Cantor and multiplication by 3. Cantor sets, which were considered an oddity when they were “discovered” at the end of the XIX century (and indeed provide the scenery to many counterexamples in analysis), are actually easily observed in dynamical systems. Consider the expanding map $E_3 : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ on the circle, which reads

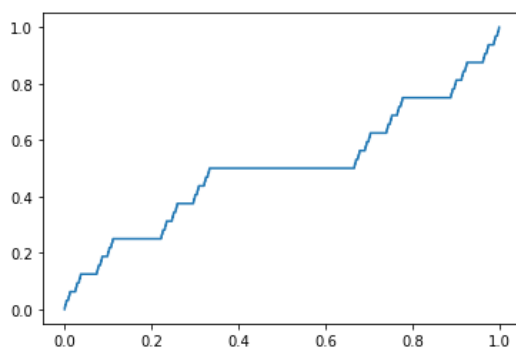
$$E_3(0.x_1x_2x_3 \cdots + \mathbb{Z}) = 0.x_2x_3x_4 \cdots + \mathbb{Z}$$

if we write points in base 3 according to $0.x_1x_2x_3 \cdots = \sum_{k=1}^{\infty} x_k/3^k$ with $x_k \in \{0, 1, 2\}$. Then the middle-third Cantor set K (once identified 0 and 1, of course) is a closed +invariant subset, since $E_3(K) = K$. It is also clear that the restriction $E_3|_K : K \rightarrow K$ is topologically mixing and also admits a dense set of periodic points. Indeed, such restriction is, almost tautologically, topologically conjugated to the shift on an alphabet of 2 letters (the letters 0 and 2 in the representation of numbers in base 3). In particular, there are many points in the circle such that the closure of their orbits under E_3 is K .

Cantor function and devil’s staircase. The continuous function $\psi : K \rightarrow [0, 1]$ defined in (9.3) is clearly non-decreasing. It can be extended to a continuous non-decreasing function $\kappa : [0, 1] \rightarrow [0, 1]$ from the unit interval onto itself, declaring that its values for $y \notin K$ are

$$\kappa(y) = \sup_{K \ni x < y} \psi(x).$$

It is called *Cantor function*⁴⁷, and its graph is known as *devil’s staircase*.



The Cantor function κ is constant on the missing intervals of the Cantor construction, i.e. in a subset of total Lebesgue measure. In particular, it has zero derivative almost everywhere, while still “growing” from 0 to 1 ! One can also show that it is uniformly continuous (actually Hölder of exponent $\log 2 / \log 3$), but not absolutely continuous.

ex: What is the length of the devil’s staircase, the graph of the Cantor function?

⁴⁷G. Cantor, De la puissance des ensembles parfaits de points: Extrait d’une lettre adressée à l’éditeur, *Acta Mathematica* 4 (1884), 381-392.

Peano curves. The homeomorphism $K \simeq K \times K$ and the continuous map $\psi : K \rightarrow [0, 1]$ of the Cantor set onto the unit interval can be combined to give a continuous map of K onto the unit square $[0, 1] \times [0, 1]$. This map may be easily extended to a continuous map of the unit interval onto the unit square (for example, declaring that the images of the missing intervals $I_i = (a_i, b_i)$ in the construction of the Cantor set are segments between the images of the boundary points a_i and b_i , which belong to the Cantor set).

This gives an elegant example of a *Peano curve* (which is not Peano original construction⁴⁸), a so called “space-filling curve”. When discovered, mathematicians were quite surprised, since they show that the “dimension”, however defined in a reasonable manner, is not invariant under continuous maps!

Cantor sets from the quadratic family. Consider the quadratic family $f_\lambda : \mathbb{R} \rightarrow \mathbb{R}$ (this time defined in the whole real line), defined by $x \mapsto \lambda x(1 - x)$, where $\lambda > 0$. The trajectory of any point outside the unit interval $I = [0, 1]$ diverges. We may therefore define the set

$$\Lambda = \bigcap_{n \geq 0} f_\lambda^{-n}(I)$$

of those points with bounded orbits. If $\lambda > 4$, a picture shows that $f_\lambda^{-1}([0, 1])$ is the disjoint union of two closed not-empty intervals I_0 and I_1 contained in $[0, 1]$. If λ is sufficiently large, it is also clear that $|f'_\lambda(x)|$ is uniformly larger than one at the points of $I_0 \cup I_1$. By induction, one can show that this implies that $f_\lambda^{-(n+1)}(I)$ is a disjoint union of 2^{n+1} compact intervals strictly contained, in pairs, in the 2^n compact intervals which form $f_\lambda^{-n}(I)$. There follows that Λ is a Cantor set, and that the restriction $f_\lambda|_\Lambda : \Lambda \rightarrow \Lambda$ is topologically conjugated to the Bernoulli shift $\sigma : \Sigma^+ \rightarrow \Sigma^+$ in the alphabet $\{0, 1\}$.

9.7 Hyperbolic automorphisms of the torus

Expanding is not necessary to produce chaos. It was Anosov, following the work by Hadamard and Hopf on the geodesic flow on surfaces with negative curvature, who discovered a large class of chaotic transformations (and flows), where chaos is due to some non-trivial way of stretching and folding.

Automorphisms of the torus. Let $\mathbb{T}^N := \mathbb{R}^N / \mathbb{Z}^N$ be the N -dimensional torus. A linear map $L : \mathbb{R}^N \rightarrow \mathbb{R}^N$ defined, in the canonical basis, by a square matrix with integer entries $A \in \text{Mat}_{N \times N}(\mathbb{Z})$, induces an endomorphism of the torus $f : \mathbb{T}^N \rightarrow \mathbb{T}^N$ according to

$$f_A(x + \mathbb{Z}^N) := Ax + \mathbb{Z}^N.$$

This is clear, since a matrix with integer entries sends the lattice \mathbb{Z}^N into itself, i.e. $A\mathbb{Z}^N \subset \mathbb{Z}^N$. If it happens that $\det A = \pm 1$, then A is invertible and its inverse A^{-1} also has integer entries. This implies that f_A is invertible too, i.e. is an *automorphism* of the torus (thought as an Abelian group).

Modular group. The existence of non-trivial automorphisms of the torus is due to arithmetical reasons. For example, orientation preserving automorphisms of the 2-dimensional torus are induced by 2×2 integer matrices with determinant one, which form the *modular group* $\text{SL}_2(\mathbb{Z})$. It is made of matrices

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

such that a, b, c, d are integers satisfying $ad - bc = 1$. But this means that rows and columns of A are made of pairs of relatively prime integers! Simple non-trivial (different from the identity) examples are

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \quad \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} \quad \begin{pmatrix} 3 & 2 \\ 1 & 1 \end{pmatrix} \quad \cdots \quad \begin{pmatrix} 5 & 7 \\ 2 & 3 \end{pmatrix} \quad \cdots$$

⁴⁸G. Peano, Sur une courbe, qui remplit toute une aire plane, *Mathematische Annalen* **36** (1890), 157-160.

And much more can be produced using the group structure. Indeed, $\mathrm{SL}_2(\mathbb{Z})$ is a “large group”, since for any pair of relative prime integers a, b the division algorithm produces an integer matrix with unit determinant and first row (a, b) . Indeed, it is one of the most interesting groups in mathematics, since it contains informations about primes, and also is related to the hyperbolic geometry of the Poincaré upper half-space \mathbb{H} . Indeed, isometries of \mathbb{H} are induced by fractional linear transformations $z \mapsto (az+b)/(cz+d)$ with $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in \mathrm{SL}_2(\mathbb{R})$, and the quotient $\mathbb{H}/\mathrm{PSL}_2(\mathbb{Z})$ is an interesting hyperbolic surface, called “modular orbifold”.

In general, orientation preserving automorphisms of the N -dimensional torus are induced by matrices $A \in \mathrm{SL}_N(\mathbb{Z})$. The homogeneous space $\mathrm{SL}_N(\mathbb{R})/\mathrm{SL}_N(\mathbb{Z})$ is the space of lattices $\Gamma \subset \mathbb{R}^N$ with unit co-volume (the volume of a fundamental region). Indeed, any $G \in \mathrm{SL}_N(\mathbb{R})$ sends the standard lattice \mathbb{Z}^N , which has a fundamental region $[0, 1]^N$ of volume one, into a lattice $G\mathbb{Z}^N$ with a fundamental region of volume one (because $\det G = 1$), and the stabilizer of the standard lattice is precisely $\mathrm{SL}_N(\mathbb{Z})$.

ex: Shows that for any pair of relatively prime integers p and q there exists a matrix $A \in \mathrm{SL}_2(\mathbb{Z})$ having p and q in the first column (or row).

Hyperbolic automorphisms of the torus. Let $f_A : \mathbb{T}^N \rightarrow \mathbb{T}^N$ be the automorphism of the torus induced by a matrix $A \in \mathrm{SL}_N(\mathbb{Z})$. If some power A^n of A has eigenvalue 1, and $v \in \mathbb{R}^N$ is a corresponding eigenvector, then the entire line $\mathbb{R}v + \mathbb{Z}^N \subset \mathbb{T}^N$ is made of periodic points (of period which divides n) of the automorphism f_A . This line may be dense in the torus or in some sub-torus, depending on the rationality properties of the coordinates of v .

A square integer matrix A , and the corresponding endomorphism of the torus, is called *hyperbolic* if it does not have eigenvalues with absolute value one, i.e. if its spectrum is disjoint from the unit circle of the complex plane. If $\det A = 1$, this also implies that (the complexification of) A has eigenvalues with both $|\lambda| > 1$ and $|\lambda| < 1$, since their product must be one. Thus, A dilates distances in some directions and contracts distances in some other directions.

Theorem 9.13. *Let $f_A : \mathbb{T}^N \rightarrow \mathbb{T}^N$ be a hyperbolic automorphism of the torus. The set of periodic points of f_A is the set $\mathbb{Q}^N/\mathbb{Z}^N$ of points with rational coordinates. In particular, $\mathrm{Per}(f_A)$ is dense in the torus.*

Proof. If $x + \mathbb{Z}^N$ is a periodic point of period $n \geq 1$, then $A^n x = x + k$, or, equivalently, $(A^n - I)x = k$, for some $k \in \mathbb{Z}^N$. If the eigenvalues of A are not roots of one, then $A^n - I$ is invertible, and it is clear that the entries of its inverse are rationals. There follows that $x = (A^n - I)^{-1}k$ has rational coordinates. Thus, periodic points are rational.

On the other side, for any fixed natural $q \geq 1$, we may consider the finite set $Q_q \subset \mathbb{T}^N$ of those points of the torus with coordinates that are integer multiples of $1/q$ (it has cardinality q^N). Since A preserves Q_q and is invertible, $f_A(Q_q) = Q_q$, i.e. f_A is a permutation of Q_q . But some power of a permutation of a finite set is the identity, so that any point of Q_q is periodic. Since the denominator q was arbitrary, this proves that all rational points in the torus are periodic. \square

Indeed, one can also compute easily the cardinalities $P_n(f_A) := \mathrm{card}(\mathrm{Per}_n(f_A))$ of n -periodic points, i.e. periodic points with period dividing n .

Theorem 9.14. *The cardinality of n -periodic point of an hyperbolic automorphism of the torus \mathbb{T}^N defined by the unimodular matrix $A \in \mathrm{SL}_N(\mathbb{Z})$ is*

$$P_n(f_A) = |\det(A^n - I)|$$

Proof. A fundamental domain for the action of the lattice \mathbb{Z}^N on the space \mathbb{R}^N is the hyper-cube $Q := [0, 1)^N$ (observe that it contains only one point with integer coordinates, the origin!). A point $x \in Q$ represents a fixed point of f_A^n if $A^n x = x + k$ with $k \in \mathbb{Z}^N$, i.e. if $(A^n - I)x \in \mathbb{Z}^N$. Thus, the number of fixed points of f_A^n is equal to the cardinality of the intersection $(A^n - I)Q \cap \mathbb{Z}^N$. Since the sets $(A^n - I)(Q)$ tile the space \mathbb{R}^N (they are rhomboids), it is clear that this cardinality is exactly the volume of $(A^n - I)(Q)$, which is equal to the absolute value of the determinant of $A^n - I$ (because Q has unit volume). \square

Arnold's cat map. The classical example of a hyperbolic automorphism of the torus \mathbb{T}^2 is induced by the unimodular matrix

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix},$$

and reads

$$f_A((x, y) + \mathbb{Z}^2) = (2x + y, x + y) + \mathbb{Z}^2.$$

It is known as Arnold's 'cat map', since Arnold, who loved cats, popularized it with the drawing of a cat in [AA67]. The eigenvalues of A are

$$\lambda_{\pm} = \frac{3 \pm \sqrt{5}}{2}.$$

Moreover, since A is symmetric, one can find eigenvectors v_{\pm} which form an orthonormal basis, and they are vectors with irrational slopes. Thus, \mathbb{R}^2 is the orthogonal direct sum $\mathbb{R}^2 = E^+ \oplus E^-$ of the eigenspaces. The linear map $x \mapsto Ax$ dilates vectors of $E^+ \setminus \{0\}$ by a factor $\lambda_+ > 1$, and contracts vectors of $E^- \setminus \{0\}$ by a factor $\lambda_- < 1$. Observe that f_A preserves areas, since $\det A = \lambda_+ \lambda_- = 1$.

Theorem 9.15. *The Arnold cat map f_A is topologically mixing, hence chaotic.*

Proof. The projections of the lines $x + E^{\pm} \subset \mathbb{R}^2$ into the torus \mathbb{T}^2 contain orbits of a minimal translation of the torus, because the slopes λ_{\pm} are irrational, and therefore they are dense in the torus. In particular, for any $\delta > 0$ there exist a length L such that any segment of length $> L$ of these lines is δ -dense in the torus (i.e. intersects any ball of radius $\geq \delta$). Let $R \subset \mathbb{R}^2$ be a small square with sides of length $\ell > 0$ parallel to the lines E^{\pm} , projecting to a small rectangle $B \subset \mathbb{T}^2$. The images $A^n(R)$ are rectangles with sides of length $\ell \cdot \lambda_+^n$ and $\ell \cdot \lambda_-^n$, still parallel to the lines E^{\pm} , respectively. If n is so large that $\ell \cdot \lambda_+^n > L$, then the complementar set $\mathbb{T}^2 \setminus f^n(R)$ does not contain balls of radius greater δ . Thus, $f^n(R)$ intersects stably any not-empty open subset of the torus. \square

ex: Use theorem 9.14 to show that the cardinality of n -periodic points of the Arnold's cat map is

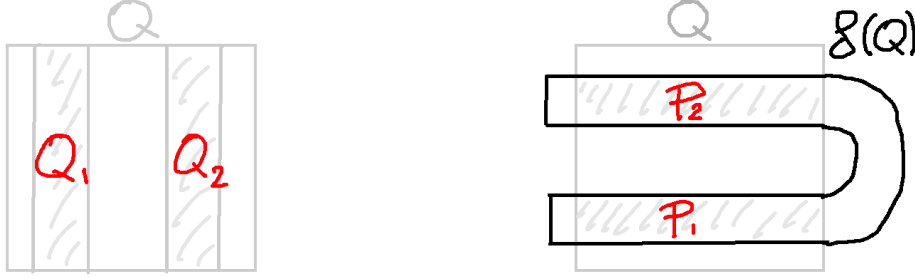
$$P_n(f_A) = \lambda_+^n + \lambda_-^n - 2$$

9.8 Horseshoes and solenoids

Here, finally, hyperbolicity and Cantor-like constructions are combined to give some of the two paradigmatic examples of chaotic dynamical systems.

Smale's horseshoe. Consider a closed rectangle Q in the euclidean plane \mathbb{R}^2 , as for example the unit square $[0, 1]^2$. We will define a map $f : Q \rightarrow \mathbb{R}^2$ which is a diffeomorphism onto its image $f(Q)$ (and which can be easily extended to a diffeomorphism of the sphere $\mathbb{S}^2 = \mathbb{R}^2 \cup \{\infty\}$ onto itself, but we'll omit this part), which exhibits a maximal invariant set $\Lambda = \bigcap_{n=-\infty}^{\infty} f^n(Q) \subset Q$ with interesting dynamics.⁴⁹

We stretch Q in the horizontal direction, squeeze Q in the vertical direction, and then bend the resulting rectangle in such a way that the intersection $Q \cap f(Q)$ is the disjoint union $P_1 \cup P_2$ of two "horizontal" rectangles P_1 and P_2 , as in the picture. Let $Q_k := f^{-1}(P_k)$ be their inverse images, so that $f^{-1}(Q) \cap Q$ is the disjoint union $Q_1 \cup Q_2$. We ask that the restrictions $f|_{Q_k} : Q_k \rightarrow P_k$ are hyperbolic affine maps, which stretch the x -direction by a factor $\alpha > 2$ and squeeze the y -direction by a factor $\beta < 1/2$. This implies that the Q_k 's are "vertical" rectangles. To achieve this, we must stretch, possibly much more, the central vertical strip, and bend the resulting rectangle as illustrated in the picture. If we want that both restrictions $f|_{Q_k}$ preserve the orientation, we must bend twice, as to form an upside-down letter G. What happens outside the vertical stripes P_k 's does not matter.



Let $\Lambda := \bigcap_{n=-\infty}^{\infty} f^n(Q)$ be the set of those points of Q with full orbit contained inside Q , and try to understand its structure.

We first observe that the intersection $f^{-1}(Q) \cap Q \cap f(Q)$ is the disjoint union of the four rectangles

$$Q_{x_{-1}x_0} = Q_{x_0} \cap P_{x_{-1}} = Q_{x_0} \cap f(Q_{x_{-1}}),$$

where the x_k 's may be 1 or 2. By induction, one easily sees that the intersection $\bigcap_{k=-n}^n f^k(Q)$ is the disjoint union of the 4^n rectangles

$$Q_{x_{-n} \dots x_{-1} x_0 x_1 \dots x_{n-1}} = \bigcap_{k=-n}^{n-1} f^{-k}(Q_{x_k}),$$

where the x_k 's belong to the alphabet $\{1, 2\}$. Each such rectangle has sides bounded above by γ^n , where $\gamma = \min\{\alpha^{-1}, \beta\} < 1/2$, and therefore area bounded above by $\gamma^{2n} < 1/4^n$. There follows that the infinite intersection Λ is a Cantor set, i.e. a compact totally disconnected perfect subset of Q , and also that it has zero area. Observe that if the P_k 's do not contain the horizontal boundaries of Q and if the Q_k 's do not contain the vertical boundaries of Q , then Λ is contained in the interior of Q .

Observe that if $x \in Q_{x_{-n} \dots x_{-1} x_0 x_1 \dots x_{n-1}}$ then $x \in Q_{x_0}$, $f(x) \in Q_{x_1}$, \dots and so on. There follows that the map $\varphi : \Sigma_2 \rightarrow \Lambda$, defined by

$$\varphi(x) = \bigcap_{k=-\infty}^{\infty} f^{-k}(Q_{x_k})$$

where $(x_k) \in \Sigma_2 = \{1, 2\}^{\mathbb{Z}}$, which is clearly a homeomorphism, is a topological conjugation between the full shift $\sigma : \Sigma_2 \rightarrow \Sigma_2$ over an alphabet of two letters and the restriction $f|_{\Lambda} : \Lambda \rightarrow \Lambda$. So,

⁴⁹S. Smale, Diffeomorphisms with many periodic points, in *Differential and Combinatorial Topology: a Symposium in Honor of Marston Morse*, Princeton Univ. Press, 1965, pp. 63-80.

Theorem 9.16. *The restriction $f|_\Lambda : \Lambda \rightarrow \Lambda$ is conjugated to the full shift over an alphabet of two letters. In particular, it is topologically mixing and admits a dense set of periodic points.*

Observe that the Smale's horseshoe is “hyperbolic”, in the sense that the differential of f at the points of Λ is a hyperbolic linear map of the plane (actually diagonal, with eigenvalues $\alpha > 1$ and $\beta < 1$).

This construction can be modified in many ways, as shown by Smale himself in [Sm67], a cornerstone in the history of dynamical systems. For example, if $f(Q) \cap Q$ is made of N horizontal rectangles, one obtains a dynamics conjugated to a shift over an alphabet of N symbols. Also the requirement that the maps between rectangles should be affine may be relaxed, still maintaining hyperbolicity, and even the rectangles may be deformed.

ex: Show that the infinite intersection $\bigcap_{k=0}^{\infty} f^k(Q)$ is the Cartesian product of the unit interval times a Cantor set $K \approx \Sigma_2$. Show that the infinite intersection $\bigcap_{k=0}^{\infty} f^{-k}(Q)$ is the Cartesian product of a Cantor set $K \approx \Sigma_2$ times the unit interval. Deduce that Λ is a Cartesian product of two Cantor sets, hence a Cantor set itself.

ex: Consider the simple case of the unit square $Q = [0, 1] \times [0, 1]$, take as horizontal rectangles $P_1 = [0, 1] \times [0, 1/3]$, and $P_2 = [0, 1] \times [2/3, 1]$ and as vertical rectangles $Q_1 := [0, 1/3] \times [0, 1]$ and $Q_2 = [2/3, 1] \times [0, 1]$. Assume that the affine maps sending the Q_k 's onto the P_k 's are

$$f|_{Q_1}(x, y) = (3x, y/3) \quad \text{and} \quad f|_{Q_2}(x, y) \mapsto (3x - 2, (y - 2)/3)$$

Show, in details, that Λ is the product $K \times K$ of two middle-third Cantor sets, and that the restriction $f|_\Lambda : \Lambda \rightarrow \Lambda$ is topologically conjugated to the full shift $\sigma : \Sigma_2 \rightarrow \Sigma_2$ over an alphabet of 2 letters.

Inverse limits/Solenoids. There is a standard construction which produces a homeomorphism out of a non-invertible continuous map (a particular case of a construction called “inverse limit”): just take the space of all possible (past, at least) histories. More precisely, consider a map $f : X \rightarrow X$, possibly not invertible, and define the space

$$X_f := \{ \dots x_{-2}x_{-1}x_0 \in X^{-\mathbb{N}_0} \text{ s.t. } f(x_n) = x_{n+1} \text{ for } n < 0 \}$$

which is a subset of the product $X^{-\mathbb{N}_0}$, equipped with the product topology. This space fibers over X , the projection $\pi : X_f \rightarrow X$ being defined by $\pi(\dots x_{-2}x_{-1}x_0) = x_0$. The fiber $\pi^{-1}(x_0)$ over a point $x_0 \in X$ is the set of possible past histories of x_0 . A map $F : X_f \rightarrow X_f$ may be defined according to

$$F(\dots x_{-2}x_{-1}x_0) := \dots x_{-2}x_{-1}x_0x_1 \quad \text{where } x_1 = f(x_0)$$

It is clear that F is invertible, its inverse being the map $\dots x_{-2}x_{-1}x_0 \mapsto \dots x_{-3}x_{-2}x_{-1}$ which forgets the first “letter”. Also, it is clear that both F and F^{-1} are continuous, so that F is a homeomorphism. The original f may be recovered just observing F at the last coordinate. Indeed, the projection π is a semi-conjugation between F and f , since by definition $F \circ \pi = \pi \circ f$. Thus, the original map f is a factor of F .

Assume that f is an expanding map of degree $\deg(f) = N > 1$ of some compact space X , so that any point has exactly N pre-images. Then the fibers $\pi^{-1}(x_0)$ are Cantor set $\{1, 2, \dots, N\}^{-\mathbb{N}} \approx \Sigma_N^+$, and may be equipped with their natural ultrametrics. Fix some $\lambda > 1$, for example $\lambda = N$. The distance between two past histories of x_0 is $d(\dots x_{-2}x_{-1}x_0 \dots y_{-2}y_{-1}x_0) = \lambda^{-n}$ if $x_{-k} = y_{-k}$ for all $k < n$ and $x_{-n} \neq y_{-n}$. The homeomorphism $F : X_f \rightarrow X_f$ is then “hyperbolic” (although we do not give a precise meaning to this word, yet), in the sense that it expands in the direction of the base X and contracts, by a factor λ^{-1} , in the direction of the fibers. Thus, this construction produces a hyperbolic homeomorphism out of an expanding map, which is a factor of the former.

When X is a circle, the resulting space is called “solenoid”. The basic example is the following.

Dyadic solenoid. Consider the simplest expanding map, the linear expanding map of the circle $E_2 : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$. The above prescription produces the *dyadic solenoid*

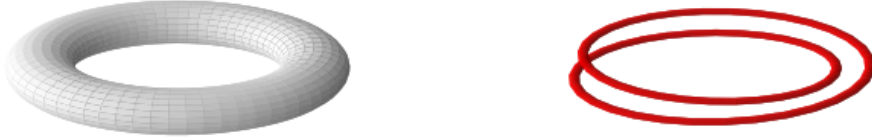
$$\mathbb{T}_2 := (\mathbb{R}/\mathbb{Z})_{E_2} = \{ \dots x_{-2}x_{-1}x_0 \in \mathbb{R}/\mathbb{Z}^{-\mathbb{N}_0} \text{ s.t. } 2x_n = x_{n+1} \pmod{1}, \text{ for } n < 0 \}$$

It fibers over the circle with fibers which are Cantor sets Σ_2^+ . The solenoid map is the homeomorphism $F_2 : \mathbb{T}_2 \rightarrow \mathbb{T}_2$ sending $\dots x_{-2}x_{-1}x_0$ to $\dots x_{-2}x_{-1}x_0(2x_0)$ (where coordinates are intended modulo 1). One easily sees that the dyadic solenoid $F_2 : \mathbb{T}_2 \rightarrow \mathbb{T}_2$ is a factor of the full shift $\sigma : \Sigma_2 \rightarrow \Sigma_2$.

Smale's solenoid. It was Smale who produced a “physical” model for the solenoid [Sm67]. Consider the solid torus $X = \mathbb{S} \times \overline{\mathbb{D}}$, parametrized by θ, z , where $\theta \in \mathbb{S} = \mathbb{R}/2\pi\mathbb{Z}$ and $z \in \overline{\mathbb{D}} = \{z \in \mathbb{C} \text{ s.t. } |z| \leq 1\}$. Consider the map $f : X \rightarrow X$ defined by

$$(\theta, z) \mapsto (2\theta, \tfrac{1}{4}z + \tfrac{1}{2}e^{i\theta})$$

(the choice of $1/4$ and $1/2$ is arbitrary, any smaller values also work). Thus, the torus is stretched by a factor two (in the circle direction), squeezed by a factor one-fourth (in the disk directions), and wrapped twice inside itself with some twisting provided by the exponential of θ , which prevents self-intersections. Indeed, it is easy to see that f defines an injection of the solid torus into its interior, $f(X)$ is a thinner torus winding twice inside X , and that the intersection between any section $D = \{\theta_0\} \times \overline{\mathbb{D}} \approx \overline{\mathbb{D}}$ of the torus and the image $f(X)$ is made of two disjoint closed disks D_{\pm} of radius $1/4$ around the points $\pm \frac{1}{2}e^{i\theta_0/2}$ of $\overline{\mathbb{D}}$ (these disks are disjoint precisely because the distance between the centers, which is one, is greater than the sum of the radii, which is one-half).



Solid torus and its image under the Smale's map.

The image $f^2(X)$, which is contained in $f(X)$, intersects each of the disks D_{\pm} in two disjoint smaller disks of radius $1/4^2$. And so on, \dots the image $f^n(X)$ is a thin solid torus winding 2^n times around X , and its intersection with a section D is made of 2^n small disjoint disks of radius $1/4^n$, two inside each of the 2^{n-1} disks which form the intersection of D with $f^{n-1}(X)$.

The *Smale attractor* is the maximal invariant set

$$\Lambda := \bigcap_{n=0}^{\infty} f^n(X).$$

It is an attractor because it admits a neighbourhood U (for example, the interior of the solid torus) such that $\bigcap_{n=0}^{\infty} f^n(U) = \Lambda$. It fibers over the circle, the projection $\pi : \Lambda \rightarrow \mathbb{S}$ being defined by $(\theta, z) \mapsto \theta$. Fibers $\pi^{-1}(\theta_0)$ are Cantor sets, naturally homeomorphic to Σ_2^+ . Indeed, one easily sees that the projection $\pi : \Lambda \rightarrow \mathbb{S}$ is a semi-conjugation between the restriction $f|_{\Lambda}$ and the expanding map E_2 .

Actually, one can show that

Theorem 9.17. *The restriction $f|_{\Lambda} : \Lambda \rightarrow \Lambda$ is topologically conjugated to the solenoid map $F_2 : \mathbb{T}_2 \rightarrow \mathbb{T}_2$ on the dyadic solenoid, and therefore it is a factor of the full shift $\sigma : \Sigma_2 \rightarrow \Sigma_2$. In particular, $f|_{\Lambda}$ is topologically mixing and admits a dense set of periodic points.*

10 Topological entropy and zeta function

10.1 Topological entropy

Entropy of coverings. Let X be a compact metric (or Hausdorff topological) space. Given an open cover \mathcal{U} , we define its *entropy*

$$H(\mathcal{U}) := \log N(\mathcal{U})$$

where $N(\mathcal{U})$ is the minimal cardinality of a subcover of \mathcal{U} , which is finite by compactness. In information theory, entropy is usually computed using base 2 logarithms, thus measured in “bits”. Here we use natural logarithms, and therefore measure entropies in “nats”.

Given two (finite) open covers \mathcal{U} and \mathcal{V} , we define their join as the open cover $\mathcal{U} \vee \mathcal{V}$ given by the opens sets $U \cap V$ with $U \in \mathcal{U}$ and $V \in \mathcal{V}$. The cardinality of such intersections that are not empty is bounded by the product of the cardinalities of the two open covers. There follows that the entropy of open covers is subadditive, i.e.

$$H(\mathcal{U} \vee \mathcal{V}) \leq H(\mathcal{U}) + H(\mathcal{V}). \quad (10.1)$$

An open cover \mathcal{V} is a *refinement* of an open cover \mathcal{U} , the notation being $\mathcal{U} \preceq \mathcal{V}$, if each $V \in \mathcal{V}$ is contained in some $U \in \mathcal{U}$. It is clear that the entropy is monotone, namely if $\mathcal{U} \preceq \mathcal{V}$ then

$$H(\mathcal{U}) \leq H(\mathcal{V}), \quad (10.2)$$

since any finite subcover of \mathcal{V} refines some subcover of \mathcal{U} of not greater cardinality.

Topological entropy. We now introduce dynamics. Let $f : X \rightarrow X$ be a continuous transformation of a compact Hausdorff space. Given an open cover \mathcal{U} , we can define the open covers $f^{-n}\mathcal{U}$, made of the open sets $f^{-n}(U)$ where $U \in \mathcal{U}$. Since inverse images of subcovers are also subcovers, it is clear that

$$H(f^{-n}\mathcal{U}) \leq H(\mathcal{U}). \quad (10.3)$$

Consider the sequence of numbers

$$h_n = H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}),$$

The sequence h_n is subadditive, i.e. $h_{n+m} \leq h_n + h_m$. Indeed, using the subadditivity (10.1) and the monotonicity (10.3) under inverse images,

$$\begin{aligned} h_{n+m} &= H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n+m-1)}\mathcal{U}) \\ &\leq H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}) + H(f^{-n}\mathcal{U} \vee f^{-n+1}\mathcal{U} \vee \dots \vee f^{-(n+m-1)}\mathcal{U}) \\ &\leq H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}) + H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(m-1)}\mathcal{U}) \\ &= h_n + h_m \end{aligned}$$

There follows from theorem 8.11 that there exists the limit

$$h(\mathcal{U}, f) := \lim_{n \rightarrow \infty} \frac{1}{n} H(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U})$$

called *entropy* of \mathcal{U} w.r.t. the map f . Assume that \mathcal{U} is finite and formed by the open sets U_1, U_2, \dots, U_N . To any point $x \in X$ and any time n we may associate a word $x_1 x_2 \dots x_n$ in the alphabet $\{1, 2, \dots, N\}$, defined according to $f^{k-1}(x) \in U_{x_k}$. Then h_n is the logarithm of the minimal number of words of length n in those letters which are necessary to describe the possible different itineraries of points of X up to time $n-1$. The limit $h(\mathcal{U}, f)$ is therefore the asymptotic exponential growth rate of those cardinalities.

The trivial cover formed by X itself has zero entropy. To get something interesting and independent of the cover we are forced to take the supremum over all covers. The *topological entropy*

of the transformation $f : X \rightarrow X$ is finally defined, according to Adler, Konheim and McAndrew⁵⁰, as

$$h_{\text{top}}(f) := \sup_{\mathcal{U}} h(\mathcal{U}, f) \quad (10.4)$$

where the sup is taken over all open covers.

It is not clear how to compute the topological entropy using the above definition. We may replace the sup by a limit, as follows. Consider a sequence of open covers \mathcal{U}_n with diameters $\text{diam}(\mathcal{U}_n) := \max_{U \in \mathcal{U}_n} \text{diam}(U) \rightarrow 0$ as $n \rightarrow \infty$. Since X is compact, any (finite) cover \mathcal{V} has a Lebesgue number ℓ (a number such that any subset of diameter $< \ell$ is contained in some $V \in \mathcal{V}$), we see that any subcover of \mathcal{V} is refined by some \mathcal{U}_n of this family (for n so large that $\text{diam}(\mathcal{U}_n) < \ell$). But then monotonicity of H implies that $H(\mathcal{V}, f) \leq H(\mathcal{U}_n, f)$, so that

$$h_{\text{top}}(f) := \lim_{n \rightarrow \infty} h(\mathcal{U}_n, f).$$

10.2 Expansiveness and generators

Expansive homeomorphisms form a rich class of topological dynamical systems where computation of the topological entropy simplifies. Their definition is clearly related to the idea of coding.

Expansive maps. A continuous transformation (or homeomorphism) $f : X \rightarrow X$ of a metric space (X, d) is *(positively) expansive* if there exists a constant $\delta > 0$ such that for all distinct $x, y \in X$ there exists a time $n \geq 0$ (or $n \in \mathbb{Z}$, in the case of a homeomorphism) such that

$$d(f^n(x), f^n(y)) > \delta$$

The greatest such constant δ is sometimes called *expansive constant* of the map f .

Equivalently, $f : X \rightarrow X$ is expansive if there exists a $\delta > 0$ such that if the orbit of two points $x, y \in X$ stays at distance $d(f^n(x), f^n(y)) < \delta$ for all times n (positive for maps or both positive and negative for homeomorphisms) then the points coincide, i.e. $x = y$. In particular, expansive maps have sensitive dependence on initial conditions.

Theorem 10.1. *Let $f : X \rightarrow X$ be an expansive homeomorphism of a compact metric space. Then the sets $\text{Per}_n(f)$ of n -periodic points are finite for every n .*

Proof. Indeed, if x and y are distinct n -periodic points, then their distance must be $d(x, y) > \delta$. But a compact metric space contains only a finite number of disjoint balls of radius $\delta/2 > 0$. \square

Expanding versus expansive. It is clear that an expanding map is expansive. Indeed, assume that f expands by a factor $\lambda > 1$ the distances between different points which are δ -near. Then δ is an expansive constant for f , since if $d(f^n(x), f^n(y)) < \delta$ for all $n \geq 0$ then the distance between x and y satisfies $\lambda^n d(x, y) < \delta$ for all $n \geq 0$, which is only possible when $d(x, y) = 0$. Thus, for example, expanding linear maps of the circle or one-sided shifts are expansive maps.

Shifts and topological Markov chains are expansive. Full or one-sided shifts are also expansive. Indeed, let $\delta > 0$ denotes the minimal distance between two points with different initial letter $x_0 \neq y_0$. If $d(\sigma^n(x), \sigma^n(y)) < \delta$ for all $n \in \mathbb{Z}$ then the two points have same letters $x_n = y_n$ for all n , hence coincide. The same holds, by the same reasoning, for topological Markov chains.

⁵⁰R.L. Adler, A.G. Konheim and M.H. McAndrew, Topological Entropy, *Transactions of the American Mathematical Society* **114** (1965), 309.

Hyperbolic automorphisms of the torus are expansive. Hyperbolic automorphisms of the torus $f_A : \mathbb{T}^2 \rightarrow \mathbb{T}^2$ are also expansive. This happens because f_A and f_A^{-1} expands by some factor $\lambda > 1$ points which are sufficiently near, say at distance $d(x, y) < \delta$, in the same unstable or stable line, respectively. If x and y are two generic points at sufficiently small distance $d(x, y) < \delta'$ (depending on δ and the slopes of the eigenvectors of A), one can join them in a unique way with an unstable segment $[x, z]$ and a stable segment $[z, y]$ of length $< \delta$. But then f_A expands by a factor λ the distance between x and z , if they are different, and f_A^{-1} expands by a factor λ the distance between y and z , if they are different. Expansivity follows easily from triangular inequality.

Isometries are not expansive. On the other side, it is clear that an isometry cannot be expansive (unless the space is finite, of course). Thus, for example, rotations of the circle (or left translation on infinite groups equipped with a left-invariant metric) are not expansive.

ex: Show that if $f : X \rightarrow X$ is expansive and $Y \subset X$ is a closed invariant subset, then also the restriction $f|_Y : Y \rightarrow Y$ is expansive.

ex: Shows that there is no expansive map $f : I \rightarrow I$ defined in a compact interval $I \subset \mathbb{R}$ (observe that such map would be locally injective, hence strictly increasing or decreasing ...)

Generators. Let $f : X \rightarrow X$ be a homeomorphism of a compact metric space X . A finite open cover $\mathcal{U} = \{U_1, U_2, \dots, U_N\}$ is called a *generator* for f if for every sequence/bi-infinite word $\{n_k\} \in \{1, 2, \dots, N\}^{\mathbb{Z}}$ the intersection

$$\bigcap_{k \in \mathbb{Z}} f^{-k} \overline{U_{n_k}}$$

contains at most one point. This means that points of x are uniquely determined by their “itinerary”, the sequences of opens sets of the cover that they visit along their history (which are not unique, if the U_n ’s overlap!).

The closures in the above definition may be omitted. Indeed, if f admits a generator \mathcal{U} , with Lebesgue number ℓ , and $\mathcal{V} = \{V_1, V_2, \dots, V_{N'}\}$ is any finite open cover with diameter $\text{diam}(\mathcal{V}) < \ell$, then clearly also the intersections $\bigcap_{k \in \mathbb{Z}} f^{-k} V_{n_k}$ contain at most one point, since each V_k is contained in the closure of some U_n .

As the word itself suggest, a “generator” allows to recover the topology of X under the iterates of the map. More precisely, given a generator $\mathcal{U} = \{U_1, U_2, \dots, U_N\}$ of cardinality N for f , we may consider the sequence of finite open covers

$$\bigvee_{k=-n}^n f^{-k} \mathcal{U} = f^n \mathcal{U} \vee \dots \vee f \mathcal{U} \vee \mathcal{U} \vee f^{-1} \mathcal{U} \vee \dots \vee f^{-n} \mathcal{U}$$

made of intersections $C_\alpha := \bigcap_{k=-N}^N f^{-k} U_{a_k}$ of finite numbers of inverse images of the U_k ’s, where $\alpha = a_{-n} \dots a_{-1} a_0 a_1 \dots a_n$ ranges between the space of words of length $2n+1$ in the letters of the alphabet $\{1, 2, \dots, N\}$. The diameters of these covers shrink to zero uniformly with n . More precisely,

Theorem 10.2. *Let \mathcal{U} be a generator for the homeomorphism $f : X \rightarrow X$ of the compact metric space X .*

i) For any $\varepsilon > 0$ there exist a time $\bar{n} \geq 0$ so large that if $n \geq \bar{n}$ and $C_\alpha \in \bigvee_{k=-n}^n f^{-k} \mathcal{U}$ then $\text{diam}(C_\alpha) < \varepsilon$.

ii) Vice-versa, for all $n \geq 0$ there exists a $\varepsilon > 0$ so small such that if $d(x, y) < \varepsilon$ then both x and y belong to some $C_\alpha \in \bigvee_{k=-n}^n f^{-k} \mathcal{U}$.

Proof. (from [Wa82] theorem 5.21) If i) is false, there exists an $\varepsilon > 0$ and two sequences of points $x_n, y_n \in \bigcap_{k=-n}^n f^{-k} U_{k,n}$, with $U_{k,n} \in \mathcal{U}$ and $n \rightarrow \infty$, at distance $d(x_n, y_n) \geq \varepsilon$. By compactness, passing to some subsequence, we can assume that both sequences converge, say $x_{n_i} \rightarrow x$ and $y_{n_i} \rightarrow y$ as $i \rightarrow \infty$, to different limits $x \neq y$, since $d(x, y) \geq \varepsilon$. Since the elements of \mathcal{U} are finite, for any fixed k the open sets U_{k,n_i} coincide with some $U_k \in \mathcal{U}$ for infinitely many n_i ’s. But then both x and y belong to $\bigcap_{k=-\infty}^{\infty} f^{-k} \overline{U_k}$, and this contradicts the fact that \mathcal{U} is a generator.

To prove ii), let δ be a Lebesgue number for the cover \mathcal{U} . By the uniform continuity of f and its inverse f^{-1} , for all n there exists an $\varepsilon > 0$ so small that if $d(x, y) < \varepsilon$ then $d(f^k x, f^k y) < \delta$ for all times $|k| \leq n$. This implies that, for all such k , both x and y belong to same $f^{-k}U_k$ for $U_k \in \mathcal{U}$, and therefore that both belong to same $\cap_{k=-n}^n f^{-k}U_k \in \bigvee_{k=-n}^n f^{-k}\mathcal{U}$. \square

Of course, given a map $f : X \rightarrow X$, it is not clear neither true that generators always exist. Indeed, it happens that existence of generators is equivalent to expansiveness⁵¹.

Theorem 10.3. *Let $f : X \rightarrow X$ be a homeomorphism of a compact metric space. Then f is expansive iff has a generator.*

Proof. Assume that f is expansive, with expansivity constant δ . Let \mathcal{U} be any finite open cover with balls of radius $\delta/2$. If both x and y belong to $\bigcap_{k \in \mathbb{Z}} f^{-k}\overline{U_{n_k}}$, then $d(f^k(x), f^k(y)) < \delta$ for every time k . By expansivity this implies that $x = y$, hence that \mathcal{U} is a generator.

Conversely, assume that \mathcal{U} is a generator for f , and let δ be its Lebesgue number. If x and y are two points such that $d(f^k(x), f^k(y)) < \delta$ for any time k , then for any k there exist a $U_{n_k} \in \mathcal{U}$ such that both $f^k(x)$ and $f^k(y)$ belong to U_{n_k} . This means that both x and y belong to $\bigcap_{k \in \mathbb{Z}} f^{-k}U_{n_k}$, and therefore that $x = y$. \square

As a consequence, expansiveness is a topological property, only the value of the expansive constant depends on the actual metric. Also, it is a property which is preserved under topological conjugacy, since a topological conjugacy sends generators to generators (if they exist).

Also interesting is that expansive homeomorphisms are factors of subspaces of full shifts ([Wa82], theorem 5.24).

ex: Show that powers f^k of an expansive homeomorphism f are also expansive (use generators).

Entropy of expansive homeomorphisms. Finally, we see that in order to compute the topological entropy of expansive homeomorphism one does not need to take a supremum: it is sufficient to consider the entropy of a generator.

Theorem 10.4. *Let $f : X \rightarrow X$ be an expansive homeomorphism of the compact metric space X . If \mathcal{U} is a generator for f , then*

$$h_{\text{top}}(f) = h(\mathcal{U}, f).$$

Proof. Let \mathcal{V} be any finite open cover, and let ℓ be its Lebesgue number. By theorem 10.2, there exists a time m so large that the diameters of the $C_\alpha \in \bigvee_{k=-m}^m f^{-k}\mathcal{U}$ are smaller than ℓ . This implies that $\bigvee_{k=-m}^m f^{-k}\mathcal{U}$ is a refinement of \mathcal{V} , so that, by monotonicity (10.2),

$$h(\mathcal{V}, f) \leq h\left(\bigvee_{k=-m}^m f^{-k}\mathcal{U}, f\right)$$

The last entropy is

$$\begin{aligned} h\left(\bigvee_{k=-n}^n f^{-k}\mathcal{U}, f\right) &= \lim_{n \rightarrow \infty} \frac{1}{n} H\left(\bigvee_{i=0}^{n-1} f^{-i}\left(\bigvee_{k=-m}^m f^{-k}\mathcal{U}\right)\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} H\left(\bigvee_{k=-m}^{m+n-1} f^{-k}\mathcal{U}\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} H\left(\bigvee_{k=0}^{2m+n-1} f^{-k}\mathcal{U}\right) \\ &= \lim_{n \rightarrow \infty} \frac{2m+n-1}{n} \frac{1}{2m+n-1} H\left(\bigvee_{k=0}^{2m+n-1} f^{-k}\mathcal{U}\right) \\ &= h(\mathcal{U}, f) \end{aligned}$$

Since $h(\mathcal{U}, f)$ is an upper bound for all the $h(\mathcal{V}, f)$'s, it is the topological entropy. \square

⁵¹H.B. Keynes and J.B. Robertson, Generators for topological entropy and expansiveness, *Math. Systems Theory* **3** (1969), 51-59.

Entropy of full shifts. Consider the full shift $\sigma : \Sigma \rightarrow \Sigma$ over an alphabet $\mathcal{A} = \{1, 2, \dots, N\}$ made of N letters. Given a letter $k \in \mathcal{A}$, let C_k be the centered cylinder made of those bi-infinite words x such that $x_0 = k$. It is clear that the finite cover $\mathcal{U} = \{C_1, C_2, \dots, C_N\}$ is a generator for σ , since the elements of $\bigvee_{k=-n}^n \sigma^{-k}\mathcal{U}$ are the centered cylinders C_α , where $\alpha = \alpha_{-n} \dots \alpha_0 \dots \alpha_n$ ranges over all the finite words of length $2n+1$ in the letters of the alphabet. Also, one easily sees that $\mathcal{U} \vee \sigma_{-1}\mathcal{U} \vee \dots \vee \sigma_{-(n-1)}\mathcal{U}$ is a cover made of N^n not-empty cylinders. There follows from theorem 10.4 that the topological entropy of the full shift over an alphabet of N letters is

$$h_{\text{top}}(\sigma) = \log N.$$

10.3 Dimensions of metric spaces

Coverings, nets and separated sets. The following notions of size of a metric space (X, d) are due to Kolmogorov's school ^{52 53}.

An ε -covering of (X, d) is a covering of $X \subset \bigcup_\alpha C_\alpha$ by subsets of diameters $\text{diam}(C_\alpha) < 2\varepsilon$. Call $C_\varepsilon(X, d)$ the minimal cardinality of an ε -covering of X .

An ε -net for (X, d) is a collection $N \subset X$ of points such that any point of X is at a distance smaller than ε from some point of N , i.e. $X \subset \bigcup_{p \in N} B_\varepsilon(p)$. Call $N_\varepsilon(X, d)$ the minimal cardinality of an ε -net for X . If X is a centered space (any subset of diameter $2r$ is contained in a ball of radius r centered in some point of X) then $N_\varepsilon(X, d) = C_\varepsilon(X, d)$.

A subset $S \subset X$ is said ε -separated (or ε -distinguishable) if its points are at a distance greater than ε from each other, i.e. if $d(p, p') > \varepsilon$ for all $p, p' \in S$ such that $p \neq p'$. The collection of disjoint balls $B_{\varepsilon/2}(p)$, where p ranges in a ε -separated set S , is also called ε -packing. Call $S_\varepsilon(X, d)$ the maximal cardinality of a set of ε -separated points inside X .

These three definitions make sense if the above extremal cardinalities are finite for every $\varepsilon > 0$, and it is not difficult to see that this happens simultaneously. The class of metric spaces with this property is called the class of *totally bounded sets* and the main examples are compact spaces.

The (base 2, for example) logarithms of these quantities have interpretations related to the probabilistic theory of transmission of signals, and are called

$$\begin{aligned} \log C_\varepsilon(X, d) & \quad \text{absolute } \varepsilon\text{-entropy of } (X, d) \\ \log N_\varepsilon(X, d) & \quad \varepsilon\text{-entropy of } (X, d) \\ \log S_\varepsilon(X, d) & \quad \varepsilon\text{-capacity of } (X, d) \end{aligned}$$

ex: Show that an ε -net defines an ε -covering, and any ε -covering determines a 2ε -net, so that

$$C_\varepsilon(X, d) \leq N_\varepsilon(X, d) \leq C_{2\varepsilon}(X, d) \quad (10.5)$$

ex: Show that a maximal ε -separated set is a ε -net, and that any ε -ball centered at a point of a minimal ε -net cannot contain more than one point of a 2ε -separated set, so that

$$S_{2\varepsilon}(X, d) \leq N_\varepsilon(X, d) \leq S_\varepsilon(X, d) \quad (10.6)$$

Box-counting dimensions. The *upper* and *lower box counting dimension* (also known as *Minkowski dimensions* or *metric dimensions*) of the metric space (X, d) are defined as

$$\overline{\dim}_b(X, d) := \limsup_{\varepsilon \searrow 0} \frac{\log N_\varepsilon(X, d)}{\log \varepsilon}$$

⁵²A. N. Kolmogorov, On certain asymptotic characteristics of completely bounded metric spaces, *Dokl. Akad. Nauk SSSR* **108**, 3 (1956), 385-389.

⁵³A.N. Kolmogorov and V.M. Tihomirov, ε -entropy and ε -capacity of sets in functional spaces, *Uspekhi Mat. Nauk* **14** (1959), 3-86. [Translated in *Amer. Math. Soc. Transl.*, series 2, **17** (1961), 277-364.]

$$\underline{\dim}_b(X, d) := \liminf_{\varepsilon \searrow 0} -\frac{\log N_\varepsilon(X, d)}{\log \varepsilon}$$

We get the same values if we substitute $S_\varepsilon(X, d)$ or $C_\varepsilon(X, d)$ to $N_\varepsilon(X)$ in the above formulas (use (10.5) and (10.6), and compare the counting functions at the values ε and 2ε).

For reasonable self-similar metric spaces the two limits coincide, and their common value $\dim_b(X)$ is simply called *box counting dimension*, and denoted by $\dim_b(X)$.

An important observation is that box dimensions do not change under scalings of the metric, i.e. if we measure distances as $d'(x, y) = \lambda d(x, y)$ instead of $d(x, y)$, for some fixed $\lambda > 0$.

ex: Show that the box-counting dimension of the n -dimensional cube $[0, 1]^n$ is what you expect, namely $\dim_b([0, 1]^n) = n$.

ex: Show that the box counting dimension of the middle-third Cantor set is $\dim_b(K) = \log 2 / \log 3$.

ex: Compute the box dimension of the space $\Sigma^+ = \mathcal{A}^{\mathbb{N}}$ of infinite words in an alphabet of N letters, equipped with the ultrametric

$$d_\lambda(x, y) = \lambda^{-\min\{k \geq 1 \text{ s.t. } x_k \neq y_k\}}.$$

Observe that a centered cylinder C_α , defined by a finite word $\alpha = \alpha_1 \alpha_2 \dots \alpha_n$ of $|\alpha| = n$ letters, is a closed ball $\overline{B}_r(x)$ of radius/diameter $r = \lambda^{-(n+1)}$ centered at any one of its points $x \in C_\alpha$, and that the distance between any two different centered cylinders C_α and C_β , defined by finite words of the same length $|\alpha| = |\beta| = n$, is $d(C_\alpha, C_\beta) \geq \lambda^{-n}$.

ex: Consider the unit interval $I = [0, 1]$ equipped with the Euclidean metric d , and define new metrics

$$d_\alpha(x, y) := d(x, y)^\alpha,$$

for $0 < \alpha \leq 1$. Verify that these are indeed metrics, and compute the box-counting dimension of the metric spaces (I, d_α) .

10.4 Topological entropy according to Bowen and Dinaburg

Bowen⁵⁴ and Dinaburg⁵⁵ adapted Kolmogorov's ideas to define an invariant of topological dynamical systems, which measures the asymptotic exponential rate of divergence of orbits. It turns out to be an alternative definition of the topological entropy.

Topological entropy according to Bowen and Dinaburg. Let $f : X \rightarrow X$ be a continuous transformation of a compact metrizable topological space X . If d is a metric on X which induces its topology, we may define a family of “dynamical metrics”, depending on time $n \geq 0$, according to

$$d_f^n(x, y) := \max_{0 \leq k \leq n} d(f^k(x), f^k(y)) \quad (10.7)$$

That is, $d_f^n(x, y)$ is the “maximal distance between the n -trajectories of x and y ”. It is clear that these metrics do not decrease with n , i.e. that $d_f^n(x, y) \leq d_f^m(x, y)$ if $n \leq m$. They are constant if f has Lipschitz constant ≤ 1 , e.g. for an isometry or a contraction, but we expect them growing with n if the transformation stretches distances in some directions.

If we fix a precision $\varepsilon > 0$, then $N_\varepsilon(X, d_f^n)$ is the “minimal number of n -orbits necessary to describe all the n -orbits with an error at most ε ”, and $S_\varepsilon(X, d_f^n)$ is the “maximal number of n -orbits which an instrument with sensibility ε can distinguish”. If X is compact, then these numbers are finite, and are monotone non-decreasing as $n \nearrow \infty$ and $\varepsilon \searrow 0$.

⁵⁴R. Bowen, Entropy for Group Endomorphisms and Homogeneous Spaces, *Transactions of the American Mathematical Society* **153** (1971), 401

⁵⁵E. Dinaburg, Relationship between topological entropy and metric entropy, *Doklady Akademii Nauk SSSR* **170** (1970), 19.

The (Bowen) *topological entropy* of the continuous transformation $f : X \rightarrow X$ of the compact metric space X is finally defined as the exponential growth rate of $N_\varepsilon(X, d_f^n)$, namely, the iterated limit

$$h_{\text{top}}^B(f) := \lim_{\varepsilon \searrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log N_\varepsilon(X, d_f^n) \quad (10.8)$$

By inequalities (10.5) and (10.6), we may replace $N_\varepsilon(X, d_f^n)$ in the above definition with $S_\varepsilon(X, d_f^n)$ or even with $C_\varepsilon(X, d_f^n)$, the cardinality of a minimal ε -cover of (X, d_f^n) . This is useful because one easily shows that

$$C_\varepsilon(X, d_f^{n+m}) \leq C_\varepsilon(X, d_f^n) \cdot C_\varepsilon(X, d_f^m) \quad (10.9)$$

Therefore, for any fixed $\varepsilon > 0$, the sequence $c_n = \log C_\varepsilon(X, d_f^n)$ is subadditive, i.e. satisfies $c_{n+m} \leq c_n + c_m$. By theorem 8.11, there exists the limit $\lim_{n \rightarrow \infty} c_n/n$. Thus, again by inequalities (10.5) and (10.6), the limsup in the definition (10.8) of topological entropy may be substituted by a limit, i.e.

$$h_{\text{top}}^B(f) = \lim_{\varepsilon \searrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log C_\varepsilon(X, d_f^n)$$

(but the corresponding limits for S_ε and N_ε need not exist!). Moreover, the different characterizations are useful to get upper bounds (from nets) and lower bounds (from separated sets) for the entropy, and therefore, in some simple cases where the two are equal, the exact value.

The notation suggests that the iterated limit which defines the topological entropy does not depend on the actual metric d , but only on the topology it induces on X . This is the case, at least for compact X 's. It is particularly important in computations, because it allows to choose a metric adapted to the transformation.

Theorem 10.5. *The Bowen topological entropy does not depend on the metric used to define the topology of the compact space X .*

Proof. Let d and d' be two equivalent metrics generating the same topology of X . Since X is compact, the identity transformation is a uniformly continuous homeomorphism between (X, d) and (X, d') . Thus, for any $\varepsilon > 0$ there exists $\delta > 0$ such that if $d'(x, y) < \delta$ then $d(x, y) < \varepsilon$. This clearly implies that $C_\varepsilon(X, d_f^n) \leq C_\delta(X, d_f'^n)$. Since the inverse homeomorphism is also uniformly continuous, the reverse inequality also holds. \square

Theorem 10.6. *If $f : X \rightarrow X$ is a factor of $g : Y \rightarrow Y$, then*

$$h_{\text{top}}^B(f) \leq h_{\text{top}}^B(g)$$

In particular, topologically conjugated dynamical systems share the same topological entropy.

Proof. This is more or less the same argument as before. Let d and d' be the metrics of X and Y , respectively. The semi-conjugation $h : Y \rightarrow X$ is uniformly continuous, because Y (and X) is compact. Therefore, for any $\varepsilon > 0$ there exists a $\delta > 0$ such that the h -image of a δ -ball in Y is contained in a ε -ball of X , i.e. $h(B_\delta(y)) \subset B_\varepsilon(h(y))$. This clearly implies that $C_\varepsilon(X, d_f^n) \leq C_\delta(Y, d_f'^n)$. Taking limits this shows the first result. The second follows changing the roles of f and g . \square

For non-compact phase spaces X , one still can define the entropy taking the supremum over compact subsets $K \subset X$. Then one shows invariance under equicontinuous homeomorphisms (see, for example, [Wa82]).

Finally, we must show that Bowen's definition recovers Adler's definition of the topological entropy, at least for continuous transformations of compact metric spaces. This is also one more proof that the Bowen topological entropy does not depend on the metric, but only on the topology.

Theorem 10.7. *The topological entropies of a continuous transformation $f : X \rightarrow X$ of a compact metric space X , defined in (10.4) and (10.8), are the same, i.e. $h_{\text{top}}(f) = h_{\text{top}}^B(f)$.*

Proof. Let \mathcal{U} is an open cover with diameter $\text{diam}(\mathcal{U}) < \varepsilon$. Two distinct points of an ε -separated set for the metric space (X, d_n) cannot belong to the same element of $\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}$. Therefore,

$$S_\varepsilon(X, d_f^n) \leq N(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}). \quad (10.10)$$

This implies that Bowen's entropy is not larger than Adler's entropy.

Conversely, let \mathcal{U} be an open cover of X , and let ℓ be its Lebesgue number. Consider any $\varepsilon < \ell$. Let x_1, x_2, \dots, x_m be a minimal ε -net for the metric space (X, d_f^n) , thus realizing $N_\varepsilon(X, d_f^n) = m$. For any such x_i , consider the open balls $B_\varepsilon(f^k(x_i))$ of radius ε centered at the images $f^k(x_i)$, for $k = 0, 1, \dots, n-1$. Any such ball is contained in some element of $\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}$, say $B_\varepsilon(f^k(x_i)) \subset U_{ik}$. The open sets

$$U_i := U_{i,0} \cap f^{-1}U_{i,1} \cap \dots \cap f^{-(n-1)}U_{i,n-1}$$

for $i = 1, 2, \dots, m$, all belong to the open cover $\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}$, and form a subcover, since the x_i 's form an ε -net. There follows that

$$N(\mathcal{U} \vee f^{-1}\mathcal{U} \vee \dots \vee f^{-(n-1)}\mathcal{U}) \leq N_\varepsilon(X, d_f^n). \quad (10.11)$$

This implies that Adler's entropy is not larger than Bowen's entropy. \square

It turns out that the second limit as $\varepsilon \rightarrow 0$, in the definition of the topological entropy, is unnecessary, provided the map is expansive and we take ε sufficiently small (half the expansive constant). We can also explicitly see this phenomenon in the simple examples below.

ex: Let $f : X \rightarrow X$ be a topological dynamical system, and d_f^n , for $n \geq 1$, be the dynamical metrics defined in (10.7). Observe that is $A \subset X$ is a set of d_f^n -diameter $< \varepsilon$ and $B \subset X$ is a set of d_f^m -diameter $< \varepsilon$, then $A \cap f^{-n}(B)$ is a set of d_f^{n+m} -diameter $< \varepsilon$. Deduce inequality (10.9).

ex: Show that $h_{\text{top}}(f^n) = n h_{\text{top}}(f)$.

Isometries have zero entropy. Contractions, isometries, or Lipschitz maps $f : X \rightarrow X$ of a metric space with Lipschitz constant ≤ 1 , have zero entropy $h_{\text{top}}(f) = 0$. This is obvious since the dynamical metrics d_n do not depend on time n .

Entropy of expanding maps of the circle. We first consider the expanding endomorphism $E_N : x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$ of the circle \mathbb{T} , with degree $N \geq 2$. One explicitly computes $N_\varepsilon(\mathbb{T}, d_{E_N}^n) \leq N^{n+m}$ and $S_\varepsilon(\mathbb{T}, d_{E_N}^n) \geq N^{n+m}$, for $\varepsilon \approx N^{-m}$. More precisely, consider the set

$$A_{n+m} = \left\{ p_k = \frac{k}{N^{n+m}} + \mathbb{Z}, \text{ with } k = 0, 1, 2, \dots, N^{n+m} - 1 \right\},$$

made of "dyadic" points of the circle with denominator N^{n+m} , of cardinality N^{n+m} . For m sufficiently large, any two successive points of A_{n+m} are at a distance $d_{E_N}^n(p_{k+1}, p_k) = N^{-m}$. Therefore, A_{n+m} is a ε -net for the metric $d_{E_N}^n$, and $N^{-m-1} < \varepsilon \leq N^{-m}$, as well as a ε -separated set for the metric $d_{E_N}^n$ and $N^{-m} < \varepsilon \leq N^{-m+1}$. These estimates imply that

$$\lim_{\varepsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{\log N_\varepsilon(\mathbb{T}, d_{E_N}^n)}{n} \leq \log N \quad \text{and} \quad \lim_{\varepsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{\log S_\varepsilon(\mathbb{T}, d_{E_N}^n)}{n} \geq \log N.$$

The two inequalities implies that $h_{\text{top}}(E_N) = \log N$. For negative N , a similar reasoning shows tha $h_{\text{top}}(E_N) = \log |N|$. Since, by theorem 9.2, a generic expanding map of circle $g : \mathbb{T} \rightarrow \mathbb{T}$ is topologically conjugated to a linear expanding map of same degree, its topological entropy is the logarithm of the absolute value of its degree, i.e.

$$h_{\text{top}}(g) = \log |\deg(g)|.$$

Entropy of Bernoulli shifts. Consider the Bernoulli shift $\sigma : \Sigma^+ \rightarrow \Sigma^+$ over an alphabet \mathcal{A} of $N \geq 2$ letters. It is convenient to use the ultrametric $d(x, y) = N^{-\min\{k \geq 1 : x_k \neq y_k\}}$ (for which cylinders are clopen balls, centered at any of their points), one explicitly computes $N_\varepsilon(\Sigma^+, d_\sigma^n) \approx N^{n+m}$ for $\varepsilon \approx N^{-m}$, since a ε -net for the metric d_σ^n on Σ^+ is given by one point for each cylinder C_α with $|\alpha| = n + m$. Therefore,

$$h_{\text{top}}(\sigma) = \log N.$$

This is as well the topological entropy of the full shift $\sigma : \Sigma_N \rightarrow \Sigma_N$.

Entropy of topological Markov chains. Now, consider the topological Markov chain $\sigma_A : \Sigma_A^+ \rightarrow \Sigma_A^+$, induced by a $N \times N$ transition matrix A . As above, for $\varepsilon \approx N^{-m}$, we must count the number of admissible cylinders C_α , i.e. defined by admissible words α , with $|\alpha| = n + m$, and this number is equal to $N_\varepsilon(\Sigma_A^+, d_{\sigma_A}^n) = \sum_{i,j} (A^{n+m-1})_{ij}$, the number of Markov paths of length $n + m$, starting with i and ending with j . The above sum is clearly bounded from above and from below by

$$c \|A^{n+m-1}\| \leq \sum_{i,j} (A^{n+m-1})_{ij} \leq C \|A^{n+m-1}\|,$$

for some positive constants c and C . The right inequality is obvious. The left inequality comes from the fact that $\sum_{i,j} (A^{n+m-1})_{ij} = \|A^{n+m-1}\|_1$, since all the entries of A are non-negative, and the fact that all norms in a finite dimensional Euclidean space are equivalent. According to Gelfand's formula, the spectral radius of a square matrix A , the maximal absolute value of its eigenvalues, can be recovered as the limit $\rho(A) = \lim_{n \rightarrow \infty} \|A^n\|^{1/n}$. Taking logarithms, we finally get

$$h_{\text{top}}(\sigma_A) = \log \rho(A).$$

ex: Compute the topological entropy of the golden ratio shift.

ex: Compute the topological entropy of the dyadic solenoid map $F_2 : \mathbb{T}_2 \rightarrow \mathbb{T}_2$ (compare with the full shift $\sigma : \Sigma_2 \rightarrow \Sigma_2$ and the expanding map $E_2 : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$).

Entropy of hyperbolic automorphisms of the torus. Consider the Arnold's cat map, the hyperbolic automorphism of the torus $f_A : \mathbb{T}^2 \rightarrow \mathbb{T}^2$ induced by the linear map of the plane defined by the matrix

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$$

To compute its entropy, we use a metric which is "adapted" to the map. Let v_\pm be the normalized eigenvectors with eigenvalues $\lambda_\pm = (3 \pm \sqrt{5})/2$. Since they form a basis (actually, an orthonormal basis), any vector of the plane is a unique superposition $x = x_+ v_+ + x_- v_-$, for some coordinates x_\pm . We define a norm in the plane as $\|x\| := \max\{|x_+|, |x_-|\}$, which is the sup norm relative to the basis v_\pm , and the corresponding metric $d(x, y) = \|x - y\|$ on \mathbb{R}^N . The ball of radius ε centered at the origin for this metric is the set of vectors x such that $|x_\pm| < \varepsilon$, a square with sides of euclidean length 2ε parallel to the eigenvectors v_\pm . The map A contracts the direction v_- by a factor λ_- , and expands the direction v_+ by a factor λ_+ . Therefore, a ball $B_\varepsilon(a; d_n)$ of radius ε and centered at the point $a = a_+ v_+ + a_- v_-$ for the dynamical metric $d_{f_A}^n$ is a rectangle with sides $|x_- - a_-| < \varepsilon$ and $|x_+ - a_+| < \varepsilon/\lambda_+^n$. The euclidean area of such a rectangle is $4\varepsilon^2/\lambda_+^n$. A covering of the torus must therefore contain at least $\lambda_+^n/4\varepsilon^2$ projections of such balls, since the torus have unit area. Thus,

$$N_\varepsilon(X, d_{f_A}^n) \geq \lambda_+^n/4\varepsilon^2.$$

In order to get upper bound, we construct explicitly a ε -net of the torus for the metric d_n . Consider, inside the unit square (which is a fundamental domain for the action of \mathbb{Z}^2 on the plane), a collection of segments parallel to v_+ at a distance ε from each other. Their cardinality is $L \leq \sqrt{2}/\varepsilon$. On each such segment, we choose points at a distance $2\varepsilon/\lambda_+^n$ from each other. Each segment contains at most $N \leq \sqrt{2} \lambda_+^n/\varepsilon$ such points (since $\sqrt{2}$ bounds the length of each segment). It is clear that this collection of points form a ε -net of the torus for the metric $d_{f_A}^n$. Thus,

$$N_\varepsilon(X, d_{f_A}^n) \leq L N = 2 \lambda_+^n/\varepsilon^2.$$

These two inequalities show that the entropy of Arnold's cat map is the logarithm of the stretching factor, namely

$$h_{\text{top}}(f_A) = \log \lambda_+.$$

A similar argument, using the Jordan normal form of the matrix A which defines an hyperbolic automorphism of the torus, shows that only the eigenvalues with absolute value larger than one count. The result is the following

Theorem 10.8. *Let $f_A : \mathbb{T}^n \rightarrow \mathbb{T}^n$ be an hyperbolic automorphism of the torus induced by the matrix $A \in \text{SL}_n(\mathbb{Z})$, with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Its topological entropy is the sum of the logarithms of the eigenvalues with absolute value greater than one, i.e.*

$$h_{\text{top}}(f_A) = \sum_k \log^+ |\lambda_k|$$

Above, we use the notation $f^+ := \max\{f, 0\}$, so that $\log^+(x) := \max\{\log x, 0\}$.

10.5 Growth of periodic orbits and zeta function

The growth of periodic orbits is also a source of invariants of a dynamical system.

Growth of periodic points. Let $f : X \rightarrow X$ be a continuous transformation of a topological space X , and let $P_n(f) := |\text{Per}_n(f)|$ be the cardinality of n -periodic points (i.e. periodic points whose period divides n). We assume implicitly that these numbers are finite for any n , hence in particular that periodic points are isolated. This is the case for expansive maps of a metrizable space, by theorem 10.1.

In typical and interesting cases, these number grow at most exponentially, and one could define an asymptotic exponential rate of growth according to

$$p(f) := \limsup_{n \rightarrow \infty} \frac{1}{n} \log^+ P_n(f)$$

You may want to compare this number, in the cases where you can compute it, with the topological entropy. Indeed, for expansive homeomorphism it is a lower bound for the topological entropy. This happens because periodic points provided separated sets.

Theorem 10.9. *Let $f : X \rightarrow X$ be an expansive homeomorphism of the compact space X . Then*

$$p(f) \leq h_{\text{top}}(f).$$

Proof. We claim that if ε is smaller than the expansive constant of f , then the set $\text{Per}_n(f)$ of n -periodic points is a ε -separated set for the dynamical metric d_n . Indeed, if x and y are points of $\text{Per}_n(f)$ at distance $d_f^n(x, y) < \varepsilon$, then also $d(f^k(x), f^k(y)) < \varepsilon$ for all times $k \in \mathbb{Z}$ (by periodicity), and therefore they must coincide by expansivity. Thus, $P_f(n) \leq S_\varepsilon(X, d_f^n)$, and the claim follows. \square

ex: Consider the expanding map of the circle $E_N : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ defined by $E_N(x + \mathbb{Z}) = Nx + \mathbb{Z}$, with $|N| \geq 2$. Verify that

$$p(E_N) = \log |N|.$$

ex: Consider the shift $\sigma : \Sigma^+ \rightarrow \Sigma^+$ over an alphabet of N letters. Verify that

$$p(\sigma) = \log N.$$

ex: Consider the Arnold's cat map $f_A : \mathbb{T}^2 \rightarrow \mathbb{T}^2$. Verify that

$$p(f_A) = \log \frac{3 + \sqrt{5}}{2}.$$

Zeta function. The collection of all these cardinalities (and not just their asymptotics) may be used to define the *Artin-Mazur zeta function*⁵⁶ of the dynamical system (motivated by the Weyl zeta function of an algebraic variety over a finite field), according to

$$\zeta_f(z) := \exp \left(\sum_{n=1}^{\infty} P_n(f) \frac{z^n}{n} \right) \quad (10.12)$$

The above formal power series defines a holomorphic function in some disk $|z| < \rho$ of the complex plane. This disk is not empty if the $P_n(f)$'s grow at most exponentially, say $P_n(f) \leq C\lambda^n$ for some constants C and $\lambda > 0$, since then it contains the disk $|z| < 1/\lambda$.

Clearly, if this convergence disk is not empty, one can recover the $P_n(f)$'s from the zeta function. It is also obvious that the zeta functions is invariant under topological conjugations, since the numbers $P_n(f)$'s are.

More interesting is that, just like Euler's product formula for the classical Riemann zeta function, also the Artin-Mazur zeta function allows a product formula. A periodic point $p \in \text{Fix}(f^n)$ defines a periodic orbit $\pi = \{p, f(p), \dots, f^{|\pi|-1}(p)\}$ of period $|\pi| := \text{card}(\pi)$ (which is the minimal $k > 0$ s.t. $f^k(p) = p$) which divides n . Thus, a periodic orbit π contains $|\pi|$ different periodic points of same period. Periodic orbits play the role of prime numbers in the product formula.

Theorem 10.10. *The Artin-Mazur zeta function is equal to the product*

$$\zeta_f(z) = \prod_{\pi} \left(1 - z^{|\pi|}\right)^{-1}$$

over all periodic orbits π .

Proof. Using the Taylor series $\log(1 - z) = \sum_{k=1}^{\infty} z^k/k$, we compute, for small $|z|$,

$$\begin{aligned} \log \prod_{\pi} \left(1 - z^{|\pi|}\right)^{-1} &= \sum_p \log \left(1 - z^{|\pi|}\right) \\ &= \sum_{\pi} \sum_{k=1}^{\infty} \frac{z^{k|\pi|}}{k} \\ &= \sum_{k=1}^{\infty} \sum_{\pi} |\pi| \frac{z^{k|\pi|}}{k|\pi|} \\ &= \sum_{n=1}^{\infty} P_n(f) \frac{z^n}{n} \end{aligned}$$

where, in the last sum, we set $n = k|\pi|$ and observed that $P_n(f)$ is the sum of the $|\pi|$'s over all periodic orbits such that $|\pi|$ divides n . \square

If the zeta function is rational, as will be the case for some simple but interesting systems, then the growth of the $P_f(n)$'s is determined by the finitely many zeros and poles of $\zeta_f(z)$.

⁵⁶M. Artin and B. Mazur, On periodic points, *Ann. of Math.* **81** (1965), 82-99.

Zeta function for the Bernoulli shift. The simplest computation is that of the shift $\sigma : \Sigma^+ \rightarrow \Sigma^+$ over an alphabet of N letters. Since $P_n(\sigma) = N^n$, we compute

$$\frac{\partial}{\partial z} \left(\sum_{n=1}^{\infty} \frac{1}{n} P_n(\sigma) z^n \right) = \frac{\partial}{\partial z} \left(\sum_{n=1}^{\infty} \frac{1}{n} (Nz)^n \right) = \sum_{n=0}^{\infty} N^{n+1} z^n = \frac{N}{1 - Nz}.$$

inside the disk $|z| < 1/N$. Integrating, we get

$$\sum_{n=1}^{\infty} \frac{1}{n} P_n(\sigma) z^n = -\log(1 - Nz)$$

where \log denotes the principal branch of the logarithm. There follows that the zeta function is

$$\zeta_{\sigma}(z) = \frac{1}{1 - Nz}.$$

ex: Show that the zeta function of the expanding map of the circle $f(x + \mathbb{Z}) = Nx + \mathbb{Z}$, where N is positive integer $N \geq 2$, is

$$\zeta_f(z) = \frac{1 - z}{1 - Nz}$$

inside the disk $|z| < 1/N$.

Zeta function for the Arnold cat map. The cardinality of n -periodic points of the Arnold's cat map $f_A : \mathbb{T}^2 \rightarrow \mathbb{T}^2$ is $P_n(f_A) = |\lambda_+^n + \lambda_-^n - 2|$, where $\lambda_{\pm} = (3 \pm \sqrt{5})/2$ are the eigenvalues of A . There follows that its zeta function is

$$\zeta_{f_A}(z) = \frac{(1 - z)^2}{(1 - \lambda_+ z)(1 - \lambda_- z)}$$

Zeta functions for topological Markov chains. Rationality of the zeta function for subshifts of finite type was proved by Bowen and Lanford⁵⁷. Here we consider the case of a topological Markov chain defined by a transition matrix A over an alphabet of N letters. Since $P_n(\sigma_A) \leq N^n$, the power series in (10.12) is certainly holomorphic in the disk $|z| < 1/N$.

Theorem 10.11 (Bowen-Lanford). *The zeta function of a topological Markov chain with transition matrix A is holomorphic in the disk $|z| < 1/\rho(A)$ and admits a meromorphic continuation to the whole complex plane, where it is given by the Bowen-Lanford formula*

$$\zeta_{\sigma_A}(z) = \frac{1}{\det(I - zA)}$$

Proof. Using theorem 9.5 and the Jordan normal form of the matrix A , one sees that

$$P_n(\sigma_A) = \text{tr}(A^n) = \sum_{k=1}^N \lambda_k^n,$$

where the λ_k 's are the eigenvalues of A , counted as many times according to their multiplicities. Arguing as above, one sees that

$$\sum_{n=1}^{\infty} \frac{1}{n} P_n(\sigma_A) z^n = \sum_{n=1}^{\infty} \sum_{k=1}^N \frac{1}{n} \lambda_k^n z^n = - \sum_{k=1}^N \log(1 - \lambda_k z)$$

⁵⁷R. Bowen and O.E. Lanford III, Zeta functions of the shift transformation, *Proc. AMS Symp. Pure Math.* **14** (1970), 43-49.

and therefore

$$\zeta_\sigma(z) = \prod_{k=1}^N \frac{1}{1 - \lambda_k z}$$

in the disk $|z| < 1/\rho(A)$, where $\rho(A) = \max_k |\lambda_k|$ denotes the spectral radius of A . Since the $(1 - \lambda_k z)$'s are the eigenvalues of $I - zA$, their product is the determinant of $I - zA$. \square

ex: Show that the Artin-Mazur zeta function of the golden ratio shift $\sigma : \Sigma_G^+ \rightarrow \Sigma_G^+$ is

$$\zeta_{\sigma_G}(z) = \frac{1}{1 - z - z^2}$$

11 Ergodicity and convergence of time means

11.1 Ergodicity

Ergodicity. Let $f : X \rightarrow X$ be an endomorphism of the measurable space (X, \mathcal{E}) . The invariant probability measure μ is said *ergodic* if any invariant event, i.e. any $A \in \mathcal{E}$ such that $f^{-1}(A) = A$, has zero or total probability, i.e. $\mu(A) = 0$ or 1 . If this happens, one also says that f is an *ergodic* endomorphism of the probability space (X, \mathcal{E}, μ) .

Equivalent conditions are also useful (see [Wa82]), which show that ergodicity is a probabilistic version of topological transitivity. Recall that an event $A \in \mathcal{E}$ of a probability space (X, \mathcal{E}, μ) is said to be *invariant mod 0* (or *quasi-invariant*) if the symmetric difference $A \Delta f^{-1}A$ has zero measure.

Theorem 11.1. *Let $f : X \rightarrow X$ be an endomorphism of the probability space (X, \mathcal{E}, μ) . The following are equivalent:*

- i) μ is ergodic.
- ii) any event $A \in \mathcal{E}$ which is invariant modulo 0 has probability $\mu(A) = 0$ or 1 .
- iii) for any $A \in \mathcal{E}$ with positive measure we have

$$\mu\left(\bigcup_{n=0}^{\infty} f^{-n}A\right) = 1,$$

- iv) for any $A, B \in \mathcal{E}$ with positive measure there exists a time $n \geq 0$ such that

$$\mu(f^{-n}A \cap B) > 0.$$

Proof. i) \Rightarrow ii) Given any event A , define

$$A_{i.o.} = \{x \in X \text{ s.t. } f^n(x) \in A \text{ infinitely often}\} = \bigcap_{n=0}^{\infty} \bigcup_{k=n}^{\infty} f^{-k}A$$

It is clear that $A_{i.o.}$ is invariant, so that, by i), it has probability zero or one. Moreover,

$$\begin{aligned} \mu(A \Delta A_{i.o.}) &= \mu(A^c \cap (\bigcap_{n=0}^{\infty} \bigcup_{k=n}^{\infty} f^{-k}A)) + \mu(A \cap (\bigcup_{n=0}^{\infty} \bigcap_{k=n}^{\infty} f^{-k}A^c)) \\ &\leq \sum_n \mu(A^c \cap f^{-n}A) + \sum_n \mu(A \cap f^{-n}A^c) \\ &\leq \sum_n \mu(A \Delta f^{-n}A) \end{aligned}$$

But quasi-invariance of A clearly implies, by induction, that all the $\mu(A \cap f^{-n}A) = 0$. Thus, $\mu(A \Delta A_{i.o.}) = 0$, so that also A has probability zero or one.

ii) \Rightarrow iii) Let $\mu(A) > 0$, and consider $A_{\infty} := \bigcup_{n=0}^{\infty} f^{-n}A$. It is clear that also $\mu(A_{\infty}) > 0$. Moreover, $f^{-1}A_{\infty} \subset A_{\infty}$ and therefore $\mu(A_{\infty} \Delta f^{-1}A_{\infty}) = 0$. By ii) we get $\mu(A_{\infty}) = 1$.

iii) \Rightarrow iv) If $\mu(A) > 0$ then, according to iii), $A_{\infty} := \bigcup_{n=0}^{\infty} f^{-n}A$ has total measure $\mu(A_{\infty}) = 1$. If also B has positive measure, then

$$0 < \mu(B) = \mu(B \cap A_{\infty}) = \mu\left(\bigcup_{n=0}^{\infty} (B \cap f^{-n}A)\right)$$

and therefore at least one of the $B \cap f^{-n}A$ has positive measure too.

iv) \Rightarrow i) Finally, consider an invariant event A with positive probability $\mu(A) > 0$. If also its complement A^c has $\mu(A^c) > 0$, by iv) we find a time $n \geq 0$ such that $0 < \mu(A^c \cap f^{-n}A) = \mu(A^c \cap A)$ (by invariance of A), an absurd. \square

Ergodicity, observables and time means. Ergodicity admits many different equivalent formulations. In particular, it says something about invariant observables and about time means.

Theorem 11.2. *Let $f : X \rightarrow X$ be an endomorphism of the measurable space (X, \mathcal{E}, μ) . The following are equivalent:*

- i) μ is ergodic.
- ii) any invariant (measurable) observable φ is constant μ -a.e.
- iii) for any observable $\varphi \in L^1(\mu)$, the time average

$$\bar{\varphi}(x) = \lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n \varphi(f^k(x))$$

exists and is equal to the mean value $\int_X \varphi d\mu$ for μ -almost any point x .

Condition iii) is the physical meaning of ergodicity, as it says that “time averages are almost everywhere constant and equal to space averages”. In particular, taking φ equal to the characteristic function of any event A , almost any trajectory spend in A a fraction of time asymptotically proportional to $\mu(A)$, as dreamed by Boltzmann in his “ergodic hypothesis”.

Conditions i) or ii) are what one usually check in order to prove ergodicity of a probability measure.

Proof. To see that iii) \Rightarrow i), let A be an invariant event, and φ its characteristic function. Invariance of A implies that φ is invariant, hence that $\bar{\varphi} = \varphi$. There follows from i) that $\mu(A) = \int_X \varphi d\mu = \varphi(x)$ for some $x \in X$, hence that $\mu(A) = 0$ or 1 , the only values of characteristic functions.

Conditions i) and ii) are clearly equivalent, since any invariant event defines an invariant function (its characteristic function), and conversely level sets of invariant functions are invariant events.

Finally, in order to show that ii) \Rightarrow iii), let $\varphi \in L^1(\mu)$ be an integrable observable. According to the Birkhoff-Khinchin ergodic theorem 7.9, the time average $\bar{\varphi}(x)$ exists for μ -almost any $x \in X$ and $\int_X \bar{\varphi} d\mu = \int_X \varphi d\mu$. Since $\bar{\varphi}$ is invariant mod 0, by iii) it is constant with probability one. This implies that $\bar{\varphi}(x) = \int_X \varphi d\mu$ for μ -almost any $x \in X$. \square

It is clear from the proof that in condition ii) we may restrict our attention to invariant observables $\varphi \in L^p(\mu)$, for any $p \geq 1$. Indeed, since the measure is finite, $L^p(\mu) \subset L^1(\mu)$ for all $p \geq 1$, and characteristic functions belong to all such spaces. This is particularly useful when we can use Fourier analysis.

Also useful is the following characterization of ergodicity as “averaged asymptotic independence”.

Theorem 11.3. *An endomorphism $f : X \rightarrow X$ of a probability space (X, \mathcal{E}, μ) is ergodic iff for any two events $A, B \in \mathcal{E}$*

$$\frac{1}{n+1} \sum_{k=0}^n \mu(f^{-k}(A) \cap B) \rightarrow \mu(A) \mu(B) \quad (11.1)$$

as $n \rightarrow \infty$.

Proof. Consider an invariant event A and take $B = A$. If (11.1) holds, then $\mu(A) = \mu(A)^2$, and therefore A has probability zero or one. Thus, f is ergodic.

Conversely, assume that f is ergodic. The Birkhoff ergodic theorem 7.9, applied to the characteristic function of A , says that

$$\lim_{n \rightarrow \infty} \frac{1}{n+1} \sum_{k=0}^n \chi_A \circ f^k \rightarrow \mu(A)$$

a.e. Multiply both sides by the characteristic function of B and integrate. Since

$$(\chi_A \circ f^k) \chi_B = \chi_{f^{-k}(A) \cap B},$$

the dominated convergence theorem implies (11.1). \square

Warning. Ergodic dynamical systems exist, and some are listed below. On the other side, to show that a physically interesting system is ergodic turns out to be extremely difficult, and very few examples are known. The most famous are some “billards”, systems of hard spheres inside a billard table interacting via elastic collisions, studied by Yakov Sinai in the sixties ...

Ergodic measures as extremal measures. We already saw that the space Prob_f of invariant probability measure is a convex and closed subset of the compact space Prob . Here, we observe that ergodic measures are the “indecomposable” elements of this set.

Theorem 11.4. *Ergodic invariant measures are the extremals of Prob_f . Namely, an invariant measure μ is ergodic iff it cannot be written as a convex combination*

$$\mu = t\mu_1 + (1-t)\mu_0$$

where $t \in (0, 1)$ of two distinct invariant measures μ_0 and μ_1

Proof. First, observe that if ν is an invariant measure which is absolutely continuous w.r.t. the ergodic measure μ , then $\nu = \mu$. Indeed, one easily verifies that the Radon-Nykodim derivative $\rho = d\nu/d\mu$ is an invariant function, and ergodicity of μ implies that it is constant and equal to one μ -a.e. Now, let μ be an ergodic measure, and assume that $\mu = t\mu_1 + (1-t)\mu_0$ for some $t \in (0, 1)$. Since both μ_0 and μ_1 are absolutely continuous w.r.t. μ , they coincide with μ , hence, are not different. To prove the converse, assume that the invariant measure μ is not ergodic, hence there exists an invariant event C such that $0 < \mu(C) < 1$. Let μ_0 and μ_1 be the “conditional probability measures” defined as $\mu_1(A) = \mu(A \cap C) / \mu(C)$ and $\mu_0(A) = \mu(A \cap C^c) / \mu(C^c)$. Clearly they are different, both are invariant, and $\mu = \mu(C)\mu_1 + (1 - \mu(C))\mu_0$. \square

Ergodic decomposition. In the first lines of the above proof, we actually showed that any two ergodic invariant measure μ and ν are either equal or “mutually singular”, namely, if $\mu \neq \nu$ then there exists a measurable set A such that $\mu(A) = \nu(A^c) = 1$ and $\mu(A^c) = \nu(A) = 0$. This suggests that maybe any invariant measure could be “disintegrated” along a partition whose atoms are the support of all the different ergodic measure, in other word that μ is a “convex combination”, namely an integral, of the ergodic measures. This is true, sometimes, but both its statement and proof are quite technical: we just quote the result.

Theorem 11.5 (Ergodic decomposition). *Let $f : X \rightarrow X$ be a continuous transformation of the compact metrizable space X . There exists a partition $\mathcal{P} = \{P_e\}_{e \in E}$ of X (modulo sets of zero measure) into invariant measurable sets indexed by a Lebesgue space E , and a measurable map $E \ni e \mapsto \mu_e \in \text{Prob}_f$ with values in the space of ergodic Borel probability measures and with the property that $\mu_e(P_e) = 1$ for any $P_e \in \mathcal{P}$, such that any invariant Borel probability measure μ can be written as an integral*

$$\mu = \int_E \mu_e d\bar{\mu}(e)$$

where $\bar{\mu}$ is some probability measure on E .

Observe that the above theorem contains the statement that any continuous transformation of a compact space admits at least one ergodic Borel probability measure.

11.2 Examples of ergodic maps

Bernoulli shift. Let $\sigma : \Sigma^+ \rightarrow \Sigma^+$ be the Bernoulli shift over the alphabet $\mathcal{A} = \{1, 2, \dots, N\}$, let $p = \{p_1, p_2, \dots, p_N\}$ be any probability on \mathcal{A} , and μ the Bernoulli invariant measure on the Borel σ -algebra \mathcal{E} defined by p .

Theorem 11.6. *The Bernoulli invariant measure μ is ergodic w.r.t. σ^+ .*

Proof. First observe that, given two centered cylinders C_α and C_β , the definition of μ implies that there exists a time $n \geq 1$ such

$$\mu(C_\alpha \cap \sigma^{-k}(C_\beta)) = \mu(C_\alpha) \cdot \mu(\sigma^{-k}(C_\beta)) = \mu(C_\alpha) \cdot \mu(C_\beta)$$

whenever $k \geq n$. Indeed, one can take $n = |\alpha| + 1$, and the above reflect the "independence" of the different trials encoded in the construction of the Bernoulli measure. By additivity, the same holds true for any couple of elements of \mathcal{E} , the algebra made of finite unions of centered cylinders. Now, assume that $A \in \mathcal{B}$ is invariant. Since any Borel set $A \in \mathcal{B}$ can be approximated in measure by an elements of \mathcal{E} , given any $\varepsilon > 0$ one can find an $A_\varepsilon \in \mathcal{E}$ such that $\mu(A \Delta A_\varepsilon) < \varepsilon$. Using the above result, we can find an $n \geq 1$ such that

$$\mu(A_\varepsilon \cap \sigma^{-n}(A_\varepsilon)) = \mu(A_\varepsilon) \cdot \mu(\sigma^{-n}(A_\varepsilon)) = \mu(A_\varepsilon)^2$$

where the last equality comes from invariance of μ . Then, observe that the symmetric difference between $A \cap \sigma^{-n}(A)$ and $A_\varepsilon \cap \sigma^{-n}(A_\varepsilon)$ is contained in $(A \Delta A_\varepsilon) \cup \sigma^{-n}(A \Delta A_\varepsilon)$. This gives

$$\begin{aligned} |\mu(A \cap \sigma^{-n}(A)) - \mu(A_\varepsilon \cap \sigma^{-n}(A_\varepsilon))| &\leq \mu(A \Delta A_\varepsilon) + \mu(\sigma^{-n}(A \Delta A_\varepsilon)) \\ &\leq 2 \cdot \mu(A \Delta A_\varepsilon) < 2\varepsilon \end{aligned}$$

which, together with

$$|\mu(A)^2 - \mu(A_\varepsilon)^2| \leq 2 \cdot \mu(A \Delta A_\varepsilon) < 2\varepsilon$$

gives

$$|\mu(A) - \mu(A)^2| < 4\varepsilon$$

Since $\varepsilon > 0$ was arbitrary, we just showed that the measure of any invariant Borel set A satisfies $\mu(A) = \mu(A)^2$, hence it is either 0 or 1. \square

Observe that this proof is very similar to the argument in the Kolmogorov zero-one law for tail events in the theory of stochastic processes.

Now, let φ_k be the characteristic function of $\{x \in \Sigma^+ \text{ s.t. } x_1 = k\}$. The observables $\varphi_k \circ \sigma^n$ form a sequence of independent and identically distributed random variables with mean p_k . One can interpret the event $\{\varphi_k \circ \sigma^n = 1\} = \{x \in \Sigma^+ \text{ s.t. } x_n = k\}$ as "success in the n -th trial", where the probability of success in each trial is p_k . The Birkhoff-Khinchin ergodic theorem, together with the ergodicity of μ , gives the result that

$$\mu \left\{ x \in \Sigma^+ \text{ s.t. } \frac{1}{n+1} (\varphi_k + \varphi_k \circ \sigma^1 + \varphi_k \circ \sigma^2 + \dots + \varphi_k \circ \sigma^n)(x) \rightarrow p_k \right\} = 1$$

which is the Kolmogorov strong law of large numbers.

Irrational translations of the torus. When dealing with maps of the torus, the most convenient tool to check ergodicity is Fourier analysis.

Theorem 11.7. *Lebesgue probability measure is ergodic for the rotation $R_\alpha : \mathbb{R}/\mathbb{Z} \rightarrow \mathbb{R}/\mathbb{Z}$ iff α is irrational.*

Proof. Let $\varphi \in L^2(\mathbb{R}/\mathbb{Z})$, and consider its Fourier series

$$\varphi(x + \mathbb{Z}) \sim \sum_{k \in \mathbb{Z}} \widehat{\varphi}(k) e^{2\pi i k x}$$

with $\sum_{n \in \mathbb{Z}} |\widehat{\varphi}(k)|^2 < \infty$. We compute

$$(\varphi \circ R_\alpha)(x + \mathbb{Z}) \sim \sum_{k \in \mathbb{Z}} \widehat{\varphi}(k) e^{2\pi i k \alpha} e^{2\pi i k x}$$

If φ is invariant then

$$\widehat{\varphi}(k) (1 - e^{2\pi i k \alpha}) = 0$$

for all $k \in \mathbb{Z}$. If α is irrational, this implies that $\widehat{\varphi}(k) = 0$ for all $k \neq 0$, and therefore that $\varphi(x + \mathbb{Z}) = \widehat{\varphi}(0)$ almost everywhere.

Conversely, when the angle is rational, say $\alpha = p/q$, one easily finds non-constant invariant observables, e.g. the character $e^{2\pi i q x}$. \square

Indeed, much more is true, as we will see in the next subsection.

It is clear that the same argument works in higher dimension. Consider a translation $R_\alpha : \mathbb{T}^n \rightarrow \mathbb{T}^n$ of the torus $\mathbb{T}^n = \mathbb{R}^n / \mathbb{Z}^n$, defined by $R_\alpha(x + \mathbb{Z}^n) = x + \alpha + \mathbb{Z}^n$, for some $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$.

Theorem 11.8. *Lebesgue probability measure is ergodic for the rotation $R_\alpha : \mathbb{R}^n / \mathbb{Z}^n \rightarrow \mathbb{R}^n / \mathbb{Z}^n$ iff $\alpha_1, \dots, \alpha_n, 1$ are rationally independent, i.e. if the scalar product $\langle k, \alpha \rangle \notin \mathbb{Z}$ for all $k \in \mathbb{Z}^n \setminus \{0\}$*

ex: Write the details of the proof of the above theorem.

Expanding endomorphisms of the circle. Let $E_N : x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$ be an expanding endomorphism of the circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ of degree N such that $|N| > 1$. Lebesgue probability measure ℓ is clearly invariant under E_N . We claim that

Theorem 11.9. *Lebesgue probability measure ℓ is an ergodic measure for E_N .*

Proof. Let $\varphi \in L^2(\mathbb{R}/\mathbb{Z})$, and consider its Fourier series

$$\varphi(x + \mathbb{Z}) \sim \sum_{k \in \mathbb{Z}} \widehat{\varphi}(k) e^{2\pi i k x}$$

with $\sum_{k \in \mathbb{Z}} |\widehat{\varphi}(k)|^2 < \infty$. We compute

$$(\varphi \circ E_N)(x + \mathbb{Z}) \sim \sum_{k \in \mathbb{Z}} \widehat{\varphi}(k) e^{2\pi i k N x}$$

If φ is invariant, then $\widehat{\varphi}(k) = \widehat{\varphi}(N^n k)$ for all $k \in \mathbb{Z}$ and all times $n \geq 0$. The Riemann-Lebesgue lemma implies that, if $k \neq 0$, then $\widehat{\varphi}(k) = \lim_{n \rightarrow \infty} \widehat{\varphi}(N^n k) = 0$. There follows that the only non-zero Fourier coefficient is $\widehat{\varphi}(0)$, so that φ is constant a.e. \square

An alternative proof, which does not use Fourier analysis, is the following.

Proof. To prove ergodicity, let A be an invariant Borel set, and assume that $\ell(A) < 1$. We must show that the complement $B = \mathbb{T} \setminus A$, that has positive measure, has indeed probability one. The argument goes as follows: if $\ell(B) > 0$, then, according to Lebesgue density theorem, B contains nearly all the mass of some nonempty interval. Namely, given any $\varepsilon > 0$, we can find an open interval I_n with length $\ell(I_n) = |N|^{-n}$ and centered at a density point of B such that

$$\ell(B \cap I_n) > (1 - \varepsilon) \cdot \ell(I_n)$$

Now observe that the restriction $E_N^n|_{I_n}$ is an injective map sending I_n onto the circle minus one point, in particular, $\ell(f^n(I_n)) = 1$. Since E_N uniformly dilatates lengths by a factor $|N|$, there follows that

$$\frac{\ell(E_N^n(B \cap I_n))}{\ell(E_N^n(I_n))} = \frac{\ell(B \cap I_n)}{\ell(I_n)}$$

Since, moreover, A is invariant, its complement B is +invariant, and this implies that the left-hand side above is equal to $\ell(B)$. There follows that

$$\ell(B) = \frac{\ell(B \cap I_n)}{\ell(I_n)} > (1 - \varepsilon)$$

and, since ε was arbitrary, that $\ell(B) = 1$. \square

Hyperbolic automorphisms of a torus. Let $f_A : \mathbb{T}^2 \rightarrow \mathbb{T}^2$ be an hyperbolic automorphism of the torus, induced by the unimodular matrix $A \in \text{SL}_2(\mathbb{Z})$. Lebesgue measure on the torus is an invariant probability measure, since $\det(A) = 1$.

Theorem 11.10. *f_A is ergodic w.r.t. Lebesgue measure iff no eigenvalue of A is a root of unity.*

Proof. Let φ be an invariant square integrable function in \mathbb{T}^2 , and consider its Fourier series

$$\varphi(x + \mathbb{Z}^n) \sim \sum_{k \in \mathbb{Z}^2} \widehat{\varphi}(k) e^{2\pi i \langle k, x \rangle}$$

One compute

$$(\varphi \circ f_A)(x + \mathbb{Z}^n) = \sum_{k \in \mathbb{Z}^2} \widehat{\varphi}(k) e^{2\pi i \langle k, Ax \rangle} = \sum_{k \in \mathbb{Z}^2} \widehat{\varphi}(kA^{-1}) e^{2\pi i \langle k, x \rangle}.$$

If φ is invariant, the Fourier coefficients must verify $\widehat{\varphi}(k) = \widehat{\varphi}(kA)$, and consequently

$$\widehat{\varphi}(k) = \widehat{\varphi}(kA^n) \quad (11.2)$$

for all times $n \geq 0$ and all wave numbers $k \in \mathbb{Z}^2$. Fix a wave number $k \in \mathbb{Z}^2$ different from 0. If the sequence of integer vectors kA^n were bounded, then some of the $v = kA^n$ would be periodic, say $vA^m = v$ (since bounded integer vectors are finite). But then A would have an eigenvalue λ such that $\lambda^m = 1$. If this is not the case, then $\|kA^n\| \rightarrow \infty$ for all $k \neq 0$. By the Riemann-Lebesgue lemma and by invariance (11.2), this implies that all the non-zero Fourier coefficients vanishes. There follows that φ is constant a.e.

Conversely, assume that A has an eigenvalue λ such that $\lambda^n = 1$, for some minimal $n \geq 1$. Then A^n has a unit eigenvalue. Since A has integer entries, the corresponding eigenvector v , satisfying $v(A^n - I) = 0$, may be taken with integer coordinates too. The trigonometric polynomial

$$\varphi(x + \mathbb{Z}^2) = \sum_{m=0}^n e^{2\pi i \langle v, A^m x \rangle}$$

is invariant, since $\langle v, A^n x \rangle = \langle vA^n, x \rangle = \langle v, x \rangle$ modulo one, and is not constant because v is different from the zero vector. Therefore, f_A is not ergodic. \square

As a consequence,

Theorem 11.11. *Hyperbolic automorphisms of the torus are ergodic w.r.t. Lebesgue measure.*

11.3 Normal numbers

Normal numbers. Lebesgue measure ℓ is ergodic w.r.t. multiplication by 10 in the unit circle, the map $E_{10}(x + \mathbb{Z}) = 10 \cdot x + \mathbb{Z}$. Identify the circle with the interval $[0, 1)$, and let $x = 0, x_1 x_2 x_3 \dots$ be the base 10 expression of a point of the circle, which is unique outside a subset of Lebesgue measure zero. For $k = 0, 1, 2, \dots, 9$, let φ_k be the characteristic function of the interval $[k/10, (k+1)/10)$, i.e. the observable which is equal to $\varphi_k(x) = 1$ if $x_1 = k$ and $\varphi_k(x) = 0$ otherwise. The time mean of φ_k is

$$\frac{1}{n+1} \sum_{j=0}^n \varphi_k(E_{10}^j(x)) = \frac{1}{n+1} \cdot \text{card} \{1 \leq j \leq n+1 \text{ s.t. } x_j = k\}$$

that is the number of k 's within the first $n+1$ digits of the decimal expansion of x . The limit as $n \rightarrow \infty$, if it exists, is the “asymptotic frequency” of k 's contained in the expansion of x . Ergodicity of μ implies that there exists a set $A_k \subset [0, 1[$ of Lebesgue measure one where the limit $\overline{\varphi_k}(x)$ exists and is equal to $\int \varphi_k d\ell = 1/10$. Since the intersection $A_0 \cap A_1 \cap \dots \cap A_9$ has still probability one, the result is that Lebesgue almost any number $x \in [0, 1)$ contains in its decimal expansion any of the letters $0, 1, 2, \dots, 9$ with asymptotic frequency $1/10$.

Actually, one could repeat the same argument considering any finite word $b = b_1 b_2 \dots b_n$ in the alphabeth $\{0, 1, 2, \dots, 9\}$, and show that there is a set $A_b \subset [0, 1[$ of probability one such that the base 10 expansion of any $x \in A_b$ contains the word b with asymptotic frequency 10^{-n} . A real number x whose base 10 expansion contains any finite word with the right asymptotic frequency is called *10-normal* (meaning “normal in base 10”). Since finite words in the alphabeth $\{0, 1, 2, \dots, 9\}$ are countable, and a countable union of zero measure sets still has zero measure, we just showed

that Lebesgue almost any real number is normal in base 10. Indeed, as first observed by Émile Borel⁵⁸,

Theorem 11.12 (Borel). *Lebesgue almost any real number is normal in every base $d \geq 2$.*

It is not so easy to give examples of normal numbers, actually of series whose sum is a normal number. Much more difficult is to show that a “given” number, such as π , $\sqrt{2}$ or e ..., is normal. Here we quote Mark Kac.⁵⁹

“As is often the case, it is much easier to prove that an overwhelming majority of objects possess a certain property than to *exhibit* even one such object. The present case is no exception. It is quite difficult to exhibit a ‘normal’ number! The simplest example is the number (written in decimal notation) $x = 0.1234567891011\dots$ where after the decimal point we write the positive integers in succession. The proof that this number is normal is by no means trivial.”

11.4 Distribution of digits in continued fractions

Continued fractions and Gauss map. Numbers in the unit interval $(0, 1]$ are uniquely represented, i.e. “coded”, by continued fractions. If we disregard rationals, which form a set of zero Lebesgue measure, we are left with infinite continued fractions $[0; a_1, a_2, a_3, \dots]$, i.e. one-sided infinite sequences $(a_1, a_2, a_3, \dots) \in \mathbb{N}^{\mathbb{N}}$. Recall that the *Gauss map* $G : (0, 1] \rightarrow [0, 1]$ is defined as

$$G(x) := 1/x - \lfloor 1/x \rfloor \quad \text{if } x \neq 0$$

(but we may also define $G(0) = 0$). Observe that for any rational $r \in \mathbb{Q}$ there exists a time n such that $G^n(r) = 0$. The infinite sequence of the continued fraction expansion of $x \sim [0; a_1, a_2, a_3, \dots] \in [0, 1] \setminus \mathbb{Q}$ is a coding of the orbit of x . Indeed,

$$G([0; a_1, a_2, a_3, \dots]) = [0; a_2, a_3, a_4, \dots]$$

This means that $a_n = \lfloor 1/G^{n-1}(x) \rfloor$, or, equivalently, $a_n = k$ if $G^{n-1}(x) \in [\frac{1}{k+1}, \frac{1}{k})$. In the language of dynamical systems, the Gauss map (restricted to the full measure set of irrationals) is conjugated to the one-sided shift $\sigma : \mathbb{N}^{\mathbb{N}} \rightarrow \mathbb{N}^{\mathbb{N}}$, the conjugation being the continued fraction representation $x \sim [0; a_1, a_2, a_3, \dots]$. In particular, the equivalence relation coming from the action of $\text{PSL}_2(\mathbb{Z})$ corresponds to “being in the same great orbit” of the Gauss map.

Ergodicity and distribution of digits. It is essentially due to Gauss himself the crucial observation that the absolutely continuous measure μ with density

$$d\mu(x) = \frac{1}{\log 2} \frac{1}{1+x} dx \tag{11.3}$$

is an invariant probability measure for G , meaning that $\mu(G^{-1}(B)) = \mu(B)$ for all Borel subsets $B \in (0, 1]$. It is sufficient to check invariance for intervals. The measure of an interval $[a, b]$ is

$$\frac{1}{\log 2} \int_a^b \frac{dx}{1+x} = \frac{1}{\log 2} \log \frac{1+b}{1+a}.$$

⁵⁸E. Borel, Les probabilités dénombrables et leurs applications arithmétiques, *Rendiconti del Circolo Matematico di Palermo* **27** (1909), 247-271.

⁵⁹M. Kac, *Statistical independence in probability, analysis, and number theory*, Carus Math. Monographs, **12**, New York 1959 (pag. 18).

The preimage $G^{-1}([a, b])$ is a union of intervals $\left[\frac{1}{b+n}, \frac{1}{a+n}\right]$ with $n = 1, 2, 3, \dots$. Therefore, its measure is

$$\begin{aligned} \frac{1}{\log 2} \sum_{n=1}^{\infty} \int_{\frac{1}{b+n}}^{\frac{1}{a+n}} \frac{dx}{1+x} &= \frac{1}{\log 2} \sum_{n=1}^{\infty} \log \frac{1 + \frac{1}{a+n}}{1 + \frac{1}{b+n}} \\ &= \frac{1}{\log 2} \sum_{n=1}^{\infty} \log(1 + a + n) - \log(a + n) + \log(b + n) - \log(1 + b + n) \\ &= \frac{1}{\log 2} \log \frac{1+b}{1+a} \end{aligned}$$

Indeed, more is true ⁶⁰

Theorem 11.13 (Knopp, 1926). *The Gauss measure μ is ergodic for the Gauss map.*

There follows from the Birkhoff-Khinchin ergodic theorem 7.9 that time-averages of integrable observables φ converge μ -a.e. and are equal to the μ -averages, i.e.

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N \varphi(G^n(x)) = \int_0^1 \varphi(x) d\mu(x) \quad \mu - \text{a.e.}$$

In particular, we may compute the distribution of digits in the continued fraction representation of a number in the unit interval, simply taking for φ the indicator function of some digit $d \in \mathbb{N}$ in $\mathbb{N}^{\mathbb{N}} \approx (0, 1]$. The result is the *Gauss-Kuzmin distribution* (conjectured by Gauss and proved by Kuzmin⁶¹, see also [Ar78]):

Theorem 11.14 (Gauss-Kuzmin, 1928). *For almost every real number x , the asymptotic frequency of the digit d in the continued fraction representation $x \sim [a_0; a_1, a_2, a_3, \dots]$ is*

$$p_d = \frac{1}{\log 2} \log \left(1 + \frac{1}{d(d+2)} \right)$$

The ergodicity of the Gauss map w.r.t. the Gauss measure imply many other “surprising” results, for clever choices of the observable φ . For example, if we choose $\varphi(x) = \log(a_1)$ then the Birkhoff averages are the geometric means of the first n partial quotients. There follows that for almost all numbers $x \sim [a_0; a_1, a_2, a_3, \dots]$ the limit $\lim_{n \rightarrow \infty} \sqrt[n]{a_1 a_2 a_3 \dots a_n}$ exists and is a constant, equal to

$$\prod_{n=1}^{\infty} \left(1 + \frac{1}{n(n+2)} \right)^{\log_2 n} \simeq 2.6854 \dots,$$

a number now called *Khinchin constant* [Kh35]. A similar result is: the n -th root of the denominators q_n of the convergents of almost all numbers converge to

$$\lim_{n \rightarrow \infty} \sqrt[n]{q_n} = e^{\pi^2/(12 \log 2)} \simeq 3.2758 \dots,$$

a number called *Khinchin-Lévy constant* ^{62 63}. On the other hand, the arithmetic mean of the partial quotients is unbounded for almost all numbers.

⁶⁰K. Knopp, Mengentheoretische Behandlung einiger Probleme der diophantischen Approximationen und der transfiniten Wahrscheinlichkeiten, *Math. Ann.* **95** (1926), 409-426.

⁶¹R.O. Kuzmin, Ob odnoi zadache Gaussa, *Doklady akad. nauk*, ser. A (1928), 375-380.

⁶²A.Y. Khinchin, Einige Sätze über Kettenbrüche, mit Anwendungen auf die Theorie der Diophantischen Approximationen, *Math. Ann.* **92** (1924), 115-125.

⁶³P. Lévy, Sur les lois de probabilité dont dependent les quotients complets et incomplets d’une fraction continue, *Bull. Soc. Math.* **57** (1929), 178-194.

11.5 Unique ergodicity and equidistribution

Unique ergodicity. A homeomorphism $f : X \rightarrow X$ of a compact metric space (X, d) is *uniquely ergodic* if it admits one, and only one, invariant Borel probability measure μ . It is clear that this unique invariant measure is ergodic.

This notion is the probabilistic counterpart of minimality, and indeed both minimality and unique ergodicity are often observed simultaneously (this means that, although equivalence of the two is false, it is not easy to think at a counterexample!). Observe that we defined unique ergodicity in the context of continuous transformations. The relevance of this notion for time means is due to the following⁶⁴

Theorem 11.15 (Oxtoby). *Let $f : X \rightarrow X$ be a homeomorphism of a compact metric space X . The following statements are equivalent:*

- i) f is a uniquely ergodic,
- ii) there exists an invariant Borel probability measure μ such that, for any continuous observable φ , the time averages $\bar{\varphi}(x)$ exist and are equal to $\int_X \varphi d\mu$ for any initial condition $x \in X$.
- iii) there exists an invariant Borel probability measure μ such that, for any continuous observable φ , the convergence

$$\frac{1}{n+1} \sum_{k=0}^n \varphi(f^k(x)) \rightarrow \int_X \varphi d\mu$$

as $n \rightarrow \infty$ holds and is uniform in $x \in X$.

Proof. It is obvious that iii) \Rightarrow ii). To show that ii) \Rightarrow i), take any invariant Borel probability measure ν , and any continuous observable φ . By invariance, $\int_X \varphi d\nu = \int_X (\varphi \circ f^k) d\nu$ for any $k \geq 0$, and therefore

$$\int_X \varphi d\nu = \int_X \left(\frac{1}{n+1} \sum_{k=0}^n (\varphi \circ f^k) \right) d\nu.$$

By ii) and the dominated convergence theorem, the limit of the r.h.s. when $n \rightarrow \infty$ is $\int_X \varphi d\mu$. Since φ was arbitrary, this implies that $\nu = \mu$. Finally, we must show that i) \Rightarrow iii). If iii) is false, there exist a continuous function ψ , a $\varepsilon > 0$, a subsequence $n_i \rightarrow \infty$ and a sequence (x_i) of points $x_i \in X$ such that

$$\left| \bar{\psi}_{n_i}(x_i) - \int_X \psi d\mu \right| \geq \varepsilon$$

(recall the notation $\bar{\varphi}_n$ for the mean averages $\frac{1}{n+1} \sum_{k=0}^n (\varphi \circ f^k)$). According to the Riesz representation theorem 7.6, there exist Borel probability measures ν_i such that $\bar{\varphi}_{n_i}(x_i) = \int_X \varphi d\nu_i$ for all continuous functions φ . By compactness of the space of invariant Borel probability measures, we may assume, up to passing to a subsequence, that the ν_i 's converge to ν in the weak* topology. One easily verifies that also the limit ν is invariant under f . Finally, one checks that ν is different from μ , since

$$\left| \int_X \psi d\nu - \int_X \psi d\mu \right| = \lim_{i \rightarrow \infty} \left| \int_X \psi d\nu_i - \int_X \psi d\mu \right| = \lim_{i \rightarrow \infty} \left| \bar{\psi}_{n_i}(x_i) - \int_X \psi d\mu \right| \geq \varepsilon.$$

The existence of two distinct invariant measures contradicts i). □

Weyl equidistribution theorem. The classical example of equidistribution was discovered by Hermann Weyl⁶⁵, and refines Dirichlet and Kronecker theorems, 8.1 and 8.9, on irrational rotations of the circle.

Theorem 11.16 (Weyl, 1916). *An irrational rotation of the circle is uniquely ergodic.*

⁶⁴J.C. Oxtoby, Ergodic sets, *Bull. Amer. Math. Soc.* **58** (1952), 116-136.

⁶⁵H. Weyl, Über die Gleichverteilung von Zahlen mod. Eins, *Math. Ann.* **77** (1916), 313-352.

Proof. Let $R_\alpha : x + \mathbb{Z} \mapsto x + \alpha + \mathbb{Z}$ with $\alpha \notin \mathbb{Q}$. We must check that time means of continuous observables φ converge uniformly to the average $\int_0^1 \varphi dx$. According to Weierstrass theorem, trigonometric polynomials are dense in the space of continuous functions of the circle. Trigonometric polynomials are finite superpositions of the characters $e_k(x + \mathbb{Z}) := e^{i2\pi kx}$, with $k \in \mathbb{Z}$ (the characters of the abelian group \mathbb{R}/\mathbb{Z}). Therefore, by a simple triangular argument, it suffices to check that uniform convergence of Birkhoff sums holds for any of the e_k . A computation gives, for $k \neq 0$,

$$\left| \frac{1}{n+1} \sum_{j=0}^n e_k(R_\alpha^j(x + \mathbb{Z})) \right| = \left| \frac{1}{n+1} \sum_{j=0}^n e^{i2\pi k j \alpha} \right| \leq \frac{2}{n+1} \cdot \frac{1}{|1 - e^{i2\pi k \alpha}|} \rightarrow 0$$

as $n \rightarrow \infty$, uniformly in x . On the other side, it is obvious that time averages of the constant character e_0 are constant and equal to 1. \square

The theorem owes its name to the fact that

$$\frac{1}{n+1} \sum_{j=0}^n \varphi(x + j\alpha) \rightarrow \int_0^1 \varphi dx$$

uniformly for any continuous function φ on the circle, and this is interpreted as saying that the sequence of points $\{x, x + \alpha, x + 2\alpha, x + 3\alpha, \dots\}$ is “equidistributed” w.r.t. Lebesgue measure.

We also observe that the convergence of time means also holds for Riemann integrable functions, since any such function ψ can be approximated by a couple of continuous functions $\varphi_- \leq \psi \leq \varphi_+$ such that the mean $\int(\varphi_+ - \varphi_-)dx$ is arbitrarily small.

On the other side, mean values of Lebesgue measurable functions need not converge. For example, the time mean of the characteristic function of the orbit of a point of the circle converge to one, while its mean value is clearly zero, since the orbit is countable.

Weyl’s theorem extends to higher-dimensional tori. Here we state a version for flows.

Linear flows on tori. Consider the torus $X = \mathbb{R}^N / \mathbb{Z}^N$ of dimension $n \geq 2$, and the linear flow $\phi_t : x + \mathbb{Z}^n \mapsto x + t\alpha + \mathbb{Z}^n$ defined by the differential equation

$$\dot{x} = \alpha$$

where $\alpha \in \mathbb{R}^N$. The “frequency vector” $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$ is said *non-resonant* if the scalar product $\langle k, \alpha \rangle = \sum_{j=1}^n \alpha_j k_j \neq 0$ for any $k \in \mathbb{Z}^n \setminus \{0\}$. As above, one can approximate any continuous function on the torus with trigonometric polynomials. One then checks that

$$\frac{1}{T} \int_0^T e^{i2\pi \langle k, x+t\alpha \rangle} dt = \frac{e^{i2\pi \langle k, x \rangle}}{i2\pi \langle k, \alpha \rangle} \frac{e^{i2\pi \langle k, \alpha \rangle T} - 1}{T} \rightarrow 0$$

as $T \rightarrow \infty$, for any $k \in \mathbb{Z}^n \setminus \{0\}$, while the time mean of the observable 1 is constant and equal to one. There follows that

Theorem 11.17. *A non-resonant linear flow on the torus is uniquely ergodic w.r.t. to Lebesgue measure.*

Digits of powers of two. Look at the successive powers of two, written in base ten:

2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096, 8192, 16384, 32768, 65536,
131072, 262144, 524288, 1048576, 2097152, 4194304, 8388608, 16777216, ...

The last digit recurs every four iterations. This happens also to the last two digits, although with a much larger period. The reason is quite dull, since there are a finite number of possibilities and the digits on the left do not interfere.

More interesting is to observe the first (non-zero) digit. Although the initial time serie

2, 4, 8, 1, 3, 6, 1, 2, 5, 1, 2, 4, 8, 1, 3, 6, 1, 2, 5, 1, 2, 4, 8, 1, ...

looks periodic, this is just an accident of the first few iterations. Moreover, any of the letters $k = 1, 2, \dots, 9$ will eventually appear, and with a definite asymptotic frequency. To see this, observe that the first digit of 2^n is equal to k iff

$$k \cdot 10^m \leq 2^n < (k+1) \cdot 10^m$$

for some integer $m \geq 0$, i.e. iff

$$\log_{10} k + m \leq n \log_{10} 2 < \log(k+1) + m$$

If we denote $\alpha = \log_{10} 2$, which is irrational, then the above inequality means that the image $R_\alpha^n(0 + \mathbb{Z})$ of the origin under the n -th iterate of the irrational rotation R_α belongs to the interval $I_k = [\log_{10} k, \log_{10}(k+1))$ of the unit circle \mathbb{R}/\mathbb{Z} . The length of this interval is

$$|I_k| = \log_{10} \left(1 + \frac{1}{k} \right).$$

By Weyl theorem 11.16, if $C_k(N)$ counts the number of times that the letters k appears as the first digit of 2^n for $1 \leq n \leq N$, then

$$\lim_{N \rightarrow \infty} \frac{C_k(N)}{N} \rightarrow \log_{10} \left(1 + \frac{1}{k} \right)$$

as $n \rightarrow \infty$. Thus, for example, the letter 1 appears about 30% of times, while the letter 7 appears only less than 6% of times.

11.6 Mixing

Finally, we describe the measure theoretical notion of mixing.

Mixing. An endomorphism $f : X \rightarrow X$ of a probability space (X, \mathcal{E}, μ) is called *mixing* if

$$\lim_{n \rightarrow \infty} \mu(f^{-n}(A) \cap B) \rightarrow \mu(A) \mu(B)$$

for all events $A, B \in \mathcal{E}$. This means that the past of any event becomes “asymptotically independent” from any other event, in the sense of probability. If we divide by $\mu(A)$ (assumed different from zero) and use invariance of the measure, this also says that

$$\frac{\mu(f^{-n}(A) \cap B)}{\mu(f^{-n}(A))} \rightarrow \mu(B)$$

as $n \rightarrow \infty$, which means that $f^{-n}(A)$ becomes uniformly distributed for large times n .

There follows from theorem 11.1 that mixing implies ergodicity.

As happens for ergodicity, this condition may be checked using integrable observables. We denote $\langle \varphi, \psi \rangle := \int_X \varphi \psi d\mu$ the L^2 inner product in $L^2(\mu)$, so that $\langle \varphi, 1 \rangle = \int_X \varphi d\mu$ is the mean value of φ . Also, we denote by U_f the isometry $\varphi \mapsto \varphi \circ f$.

Theorem 11.18. *Let $f : X \rightarrow X$ be an endomorphism of the measurable space (X, \mathcal{E}, μ) . The following are equivalent:*

- i) f is mixing
- ii) for any observables $\varphi, \psi \in L^2(\mu)$,

$$\langle U_f^n \varphi, \psi \rangle \rightarrow \langle \varphi, 1 \rangle \langle 1, \psi \rangle \quad (11.4)$$

as $n \rightarrow \infty$,

Condition ii) is easier to check using Fourier analysis (when possible). The implication ii) \Rightarrow i) is obvious taking characteristic functions, which are square integrable. The other is also true (a proof can be found in [Wa82]), but we don't need it.

Expanding endomorphisms of the circle. Let $E_N : x + \mathbb{Z} \mapsto Nx + \mathbb{Z}$ be an expanding endomorphism of the circle $\mathbb{T} = \mathbb{R}/\mathbb{Z}$ of degree N such that $|N| > 1$. Lebesgue probability measure ℓ is invariant under E_N , and also ergodic. We claim that

Theorem 11.19. E_N is mixing.

Proof. It is easy to check (11.4) for characters $e_k(x + \mathbb{Z}) = e^{2\pi i k x}$'s. Indeed, the inner product

$$\langle U_f^n e_k, e_m \rangle = \int_0^1 e^{2\pi i (N^n k - m)x} dx$$

vanishes for sufficiently large n as long as k or m are different from zero, just like the product $\langle e_k, 1 \rangle \langle 1, e_m \rangle$. On the other side, it also happens that $\langle U_f^n e_0, e_0 \rangle = 1 = \langle e_0, 1 \rangle \langle 1, e_0 \rangle$.

The claim follows from linearity and from denseness of trigonometric polynomials in L^2 . \square

ex: Show that the restriction of a map f on a periodic orbit $\{p_0, p_1, \dots, p_{n-1}\}$ of period $n \geq 2$ is not mixing w.r.t. the Dirac measure $\mu = \frac{1}{n}(\delta_{p_0} + \delta_{p_1} + \dots + \delta_{p_{n-1}})$.

ex: Show that an irrational rotation of the circle is not mixing w.r.t. the Lebesgue measure (use harmonics).

References

- [AA67] V.I. Arnold and A. Avez, *Problèmes ergodiques de la mécanique classique*, Gauthier-Villars, 1967.
- [Ah78] L.V. Ahlfors, *Complex Analysis*, McGraw-Hill, 1979.
- [Ap69] T.M. Apostol, *Calculus*, John Wiley & Sons, New York 1969.
- [Ar78] V.I. Arnold, *Metodi geometrici della teoria delle equazioni differenziali ordinarie*, Editori Riuniti - MIR, Roma 1978.
- [Ar79] V.I. Arnold, *Metodi matematici della meccanica classica*, Edizioni MIR - Editori Riuniti, Roma 1979.
- [Ar85] V.I. Arnold, *Equações diferenciais ordinárias*, MIR 1985.
- [AS64] M. Abramowitz and I. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, Dover, 1964.
- [Ax97] S. Axler, *Linear Algebra Done Right*, second edition, Springer, 1997.
- [Bi65] P. Billingsley, *Ergodic Theory and Information*, Wiley, 1965.
- [BV12] L. Barreira e C. Valls, *Sistemas dinâmicos, uma introdução*, IST Press, 2012.
- [BN05] P. Buttà e P. Negrini, *Note del corso di Sistemi Dinamici*, Università di Roma “La Sapienza”, 2005.
- [BS03] M. Brin and G. Stuck, *Introduction to Dynamical Systems*, Cambridge University Press, 2003.
- [CG93] L. Carleson and T.W. Gamelin, *Complex dynamics*, UTX, Springer-Verlag, 1993.
- [Chaos] P. Cvitanović, R. Artuso, P. Dahlqvist, R. Mainieri, G. Tanner, G. Vattay, N. Whelan and A. Wirzba, *Chaos: Classical and Quantum*, <http://ChaosBook.org> (Niels Bohr Institute, Copenhagen 2008).
- [CR48] R. Courant and H. Robbins, *What is mathematics?*, Oxford University Press, 1948.
- [De89] R.L. Devaney, *An introduction to chaotic dynamical systems*, Addison-Wesley, 1989.
- [De92] R.L. Devaney, *A first course in chaotic dynamical systems*, Addison-Wesley, 1992.
- [EW10] M. Einsiedler and T. Ward, *Ergodic Theory with a view towards Number Theory*, GTM **259**, Springer, 2010.
- [Fa85] K. J. Falconer, *The geometry of fractal sets*, Cambridge University Press, 1985.
- [Fe63] R.P. Feynman, R.B. Leighton and M. Sands, *The Feynman lectures on physics*, Addison-Wesley, Reading, 1963.
- [Gh07] E. Ghys, *Résonances et petits diviseurs*, in *L'héritage scientifique de Kolmogorov*, Berlin 2007.
- [GL] E. Ghys and J. Leys, *Lorenz and Modular Flows: A Visual Introduction*
- [Ha74] P. Halmos, *Measure theory*, Springer-Verlag. New York 1974.
- [Har49] G.H. Hardy, *Divergent series*, Oxford University Press, 1949.
- [HK03] B. Hasselblatt and A. Katok, *A first course in dynamics: with a panorama of recent developments*, Cambridge University Press, 2003.
- [HS74] M.W. Hirsch and S. Smale, *Differential equations, dynamical systems and linear algebra*, Academic Press, 1974.

- [HSD04] M.W. Hirsch, S. Smale and R.L. Devaney, *Differential Equations, Dynamical Systems, and an Introduction to Chaos*, 2nd ed., Elsevier Academic Press, 2004.
- [HW59] G.H. Hardy and E.M. Wright, *An Introduction to the Theory of Numbers*, fourth edition, Oxford University Press, 1959.
- [Kh35] A.Ya. Khinchin, *Continued Fractions*, 1935 [translation by University of Chicago Press, 1954].
- [KH95] A. Katok and B. Hasselblat, *Introduction to the modern theory of dynamical systems*, Encyclopedia of mathematics and its applications, Cambridge University Press, 1995.
- [Kn05] O. Knill, *Dynamical systems*, Harvard University, 2005.
- [La87] S. Lang, *Linear Algebra*, Third Edition, UTM Springer, 1987.
- [Ma75] B. Mandelbrot, *Les object fractals: forme, hasard, et dimension*, Flammarion, Paris 1975.
- [Mat95] P. Mattila *Geometry of Sets and Measures in Euclidean Spaces: Fractals and rectifiability*, Cambridge University Press, 1995.
- [Mi91] J. Milnor, *Dynamics in one complex variable*, IMS preprint, 1991.
- [MS93] W. de Melo and W. van Strien, *One-Dimensional Dynamics*, Springer-Verlag, 1993.
- [MSW02] D. Mumford, C. Series and D. Wright, *Indra's Pearls: The Vision of Felix Klein*, Cambridge University Press, 2002.
- [PM78] J. Palis jr. e W. de Melo, *Introdução aos sistemas dinâmicos*, Projeto Euclides, IMPA, 1978.
- [Ro99] J.C. Robinson, *Dynamical Systems, Stability, Symbolic Dynamic and Chaos*, CRC Press, Cambridge 1999.
- [Ro04] J.C. Robinson, *An introduction to ordinary differential equations*, Cambridge University Press, 2004.
- [Ru87] W. Rudin, *Real and complex analysis*, McGraw-Hill, 1987.
- [Sm67] S. Smale, *Differentiable dynamical systems*, *Bull. of the AMS* **73** (1967), 747-817.
- [SS03] E.M. Stein and R. Shakarchi, *Fourier Analysis. An Introduction*, Princeton University Press, 2003.
- [To16] M.J. Torres, *Tópicos de Sistemas Dinâmicos*, DMA Publicações Pedagógicas, RepositórioUM, 2016. <http://hdl.handle.net/1822/43956>
- [St94] S.H. Strogatz, *Nonlinear Dynamics and Chaos*, Addison-Wesley, 1994.
- [Wa82] P. Walters, *An Introduction to ergodic theory*, Graduate Texts in Math. **79**, Springer-Verlag, 1982