

# Mathematical Methods in Biology and Neurobiology

Jürgen Jost<sup>1</sup>

January 23, 2007

<sup>1</sup>Max Planck Institute for Mathematics in the Sciences, Inselstr.22, 04103 Leipzig,  
Germany, jost@mis.mpg.de

Copyright by the author. No reproduction or distribution without the author's permission



# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Theses about biology . . . . .	5
1.2	Fundamental biological concepts . . . . .	6
1.3	A classification of mathematical methods . . . . .	6
<b>2</b>	<b>Discrete structures</b>	<b>9</b>
2.1	Graphs and networks . . . . .	9
2.1.1	Graphs in biology . . . . .	9
2.1.2	Definitions and qualitative properties . . . . .	10
2.1.3	The graph Laplacian and its spectrum . . . . .	13
2.2	Descendence relations . . . . .	23
2.2.1	Trees and phylogenies . . . . .	23
2.2.2	Genealogies (pedigrees) . . . . .	34
2.2.3	Gene genealogies (coalescents) . . . . .	35
<b>3</b>	<b>Stochastic processes</b>	<b>41</b>
3.1	Random variables . . . . .	41
3.2	Random processes . . . . .	46
3.3	Poisson processes and neural coding . . . . .	47
3.4	Branching processes . . . . .	53
3.5	Random graphs . . . . .	59
<b>4</b>	<b>Pattern formation</b>	<b>67</b>
4.1	Partial differential equations . . . . .	67
4.2	Diffusion and random walks . . . . .	80
4.3	Dynamical systems . . . . .	89
4.3.1	Systems of ordinary differential equations . . . . .	89
4.4	Reaction-diffusion systems . . . . .	106
4.4.1	Reaction-diffusion equations . . . . .	106
4.4.2	Travelling waves . . . . .	111
4.4.3	Reaction-diffusion systems . . . . .	113
4.4.4	The Turing mechanism . . . . .	119
4.5	Continuity and Fokker-Planck equations . . . . .	126



# Chapter 1

## Introduction

### 1.1 Theses about biology

**Thesis 1.** *Biological structures are aggregate structures. Therefore, biological laws are not basic ones that do not admit exceptions, but rather emerging from some lower scale.*

**Thesis 2.** *Biological processes intertwine stochastic effects and deterministic dynamics. Randomness can support order while deterministic processes can be unpredictable, chaotic. The question then is at which level regularities emerge.*

**Thesis 3.** *Large populations of discrete units can be described by continuous models and, conversely, invariant discrete quantities can emerge from an underlying continuous substrate.*

**Thesis 4.** *Fundamental biological concepts, like fitness or information, are relative and not absolute ones.*

**Thesis 5.** *Fundamental biological quantities do not satisfy conservation laws. Those rather appear as external constraints.*

**Thesis 6.** *Biological systems interact with their environments and are thermodynamically open. Biological structures sustain the processes that reproduce them and are therefore operationally closed.*

**Thesis 7.** *Biological structures are results of historical processes. It is the task of biological theory to distinguish the regularities from the contingencies.*

**Thesis 8.** *The abstract question posed to mathematics by biology is the one of structure formation. This needs to be understood as a process because living structures are not at thermodynamic equilibrium.*

**Thesis 9.** *Gathering biological data without guiding concepts and theories is useless.*

## 1.2 Fundamental biological concepts

1. The **gene** is the unit of coding, function, and inheritance. As such, it links molecular biology and evolutionary biology. The Neodarwinian Synthesis combined Mendel and Darwin. Modern molecular biology seems to offer a more basic perspective.
2. The **cell** is the unit of metabolism. It constitutes the basic operationally closed, autopoietic system in biology. Modern biology struggles to understand the cell on the basis of their molecular constituents, DNA, RNA, and polypeptides (proteins). Multicellular organisms emerge through a partial suppression of the autonomy of the constituting cells.
3. The **species** represents the balance between the diverging effects of genetic mutations and selection at the organismic or other levels and the converging mechanism of sexual recombination. It is the arena of population biology, a child of the Neodarwinian Synthesis and the first success of mathematical models in biology. It is also important in ecology.

The **organism**, in fact, is the carrier of genes, the organisation of cells, and the member of a species. It thus links the three fundamental biological concepts. It is also a, but not the exclusive, unit of selection.

It seems that neurobiology has not yet identified such a fundamental concept, but perhaps the **spike** can be considered as the basic event of information transmission, and the **synapse** as the basic structure supporting this.

## 1.3 A classification of mathematical methods

1. Discrete structures → **Algebra**
  - (a) Static structures
    - i. Algebraic concepts: Combination and composition of objects
    - ii. Graphs and networks, including phylogenetic trees
    - iii. Information
    - iv. Discrete invariants of continuous structures and dynamical processes
  - (b) Discrete processes (Cellular automata, Boolean networks, finite state machines,...)
  - (c) Game theory as the formalisation of competition
2. Continuous methods → **Analysis**
  - (a) Deterministic dynamical processes
    - i. Continuous states enable phase transitions and bifurcations, that is, qualitative structural changes resulting from small underlying variations

- ii. Continuous states and time: Ordinary differential equations and other dynamical systems
- iii. Continuous spatial structures: Partial differential equations (example: Reaction-diffusion equations)
- (b) Stochastic analysis
  - i. Stochastic processes (while stochastic processes may also operate on discrete quantities, the concept of probability is a continuous one)
  - ii. Population processes: averaging over stochastic fluctuations in lower level dynamics
  - iii. Optimisation schemes with stochastic ingredients: Genetic and other evolutionary algorithms, swarm algorithms for distributed search, certain neural networks,...
  - iv. Statistical methods for the analysis of biological data
- 3. Hybrid models
  - (a) Difference equations (continuous states, but discrete time)
  - (b) Dynamical networks (dynamical systems coupled by a graph), in particular neural networks
- 4. System theory as a global unifying perspective?

According to the preceding list that not all mathematical subjects seem to be relevant for biology. Classical algebraic structures occur in a cursory manner at best, and one of the deepest branches, number theory and arithmetics, is entirely absent. Also, the third area of mathematics besides algebra and analysis, namely **geometry**, is entirely missing in our list. This does not mean that it is irrelevant in a similar manner as number theory is because its objects are not found in biology, but rather that it plays only a somewhat subordinate role in comparison with analysis and discrete mathematics. Organisms and their constitutive biological structures like cells are living and interacting in space, and are defining and shaping their own spaces like architectural structures which is constitutive for morphology. Symmetries and invariances, the merging ground of algebra and geometry, are important issues for the neurobiology underlying cognition. as well as for many classification purposes. One can, of course, also consider more abstract spaces of relationships, like state spaces in dynamical systems. Thus, there is some role for geometry after all.





## Chapter 2

# Discrete structures

### 2.1 Graphs and networks

#### 2.1.1 Graphs in biology

A graph is the mathematical structure representing binary relationships between discrete elements. These elements are the vertices of the graph, and the relationships are encoded as connections or edges between vertices. Such a graph can then be a network, that is, the substrate of dynamical interactions carried by the edges between processes located at the vertices. Biological applications abound.

In neural networks, the vertices stand for neurons, and the edges for synaptic connections between them. The interaction is the electrochemical transmission of pulsed dynamical activity, the spikes generated in the neurons. This activity is considered to be the carrier of information, enabling cognitive processes, but the precise identification of the information inside that dynamical activity remains unclear at present.

At smaller scales, the vertices can represent molecules like proteins, and the edges again interactions between them. The vertices can also stand for genes, and the edges for correlations in expression patterns indicating functional interactions.

At larger scales, the vertices can be the members of a population, and the edges social or other interactions, like mating. For a population with separate sexes, we then have a bipartite graph, that is, one with two distinct classes of elements such that edges exist only between members of opposite classes, but not inside one class.

At the still larger scale of ecosystems, the vertices can represent species, and the edges stand for trophic interactions. The graph then encodes a foodweb.

Another important class of biological graphs are the phylogenetic trees that turn genetic or other similarities between species into descendance relations from common ancestors. For individual descendance relations inside a sexually

recombining species we rather have pedigrees because each individual then has two parents which in turn may have more than one offspring.

### 2.1.2 Definitions and qualitative properties

We now display some formal definitions and start with the simplest situation. A graph  $\Gamma$  is a pair  $(V, E)$  of a finite set  $V$  of vertices or nodes and a set  $E$  of unordered pairs, called edges or links, of different elements of  $V$  (and we assume  $E \neq \emptyset$  to make the graph nontrivial). Thus, when there is an edge  $e = (i, j)$  for  $i, j \in V$ , we say that  $i$  and  $j$  are connected by the edge  $e$  and that they are neighbours,  $i \sim j$ . Defining edges as unordered pairs of vertices means that we consider  $(i, j)$  and  $(j, i)$  as the same pair. Thus, the neighborhood relation is symmetric. Requiring that the vertices connected by an edge be different then means that there are no edges connecting a vertex to itself. Thus, the neighborhood relation is not reflexive. In general, it is not transitive either, that is,  $i \sim j$  and  $j \sim k$  need not imply  $i \sim k$ . The degree  $n_i$  of the vertex  $i$  is the number of its neighbours. Also, the order  $|\Gamma|$  is the number of vertices in  $\Gamma$ , i.e., the cardinality of the vertex set  $V$ .

So far, we are assuming that the edges are undirected, that is, the edge  $(i, j)$  is the same as  $(j, i)$ . One may, naturally, also consider directed graphs, that is, where an edge  $e = (i, j)$  is considered to go from  $i$  to  $j$  rather than connect  $i$  and  $j$  in a symmetric manner. For example, this is appropriate for formalize neurobiological networks because synapses between neurons are directed, starting at the presynaptic neuron and going to the postsynaptic one. In addition, synapses have strengths or weights, and so, we can also consider weighted graphs where each edge  $e$  carries a weight or label  $w_e$  that indicates its strength. In fact, we may then also allow that some of the weights are negative. In a neural network, an edge with a negative weight would represent an inhibitory synapse.

Of course, every unweighted graph becomes a weighted one by assigning the weight 1 to every edge. An undirected graph with positive weights becomes a metric space by identifying each edge  $e$  with the interval of length  $(w_e)^{-1}$ . In particular, an unweighted graph then is a metric space where each edge is isometric to the unit interval. The distance between vertices then equals the length of the shortest path joining them. In particular, neighbors in the graph have distance 1.

We shall start with undirected and unweighted graphs as the simplest case. In the definition, we require that our graphs  $\Gamma$  be finite, a biologically directly plausible assumption. Moreover, we shall assume, unless stated to the contrary, that they are connected. That means that for every pair of distinct vertices  $i, j$  in  $\Gamma$ , there exists a path between them, that is, a sequence  $i = i_0, i_1, \dots, i_m = j$  of distinct vertices such that  $i_{\nu-1} \sim i_\nu$  for  $\nu = 1, \dots, m$ . Since we can decompose graphs that are not connected into their connected components, the connectivity assumption is no serious restriction.

An obvious way of representing a graph  $\Gamma$  with vertices  $i = 1, \dots, N$  is provided by its adjacency matrix  $A = (a_{ij})$ . In the unweighted case, we put  $a_{ij} = 1$  when there is an edge from  $i$  to  $j$  and  $= 0$  else. We have  $a_{ii} = 0$  because we

exclude self-loops of vertices, and  $\Gamma$  is undirected iff  $a_{ij} = a_{ji}$  for all  $i, j$ . In the weighted case, we simply put  $a_{ij} = w_{ij}$ , the weight of the edge from  $i$  to  $j$ . Of course, most large graphs arising in applications are sparse, that is, between most pairs  $i, j$ , there is no edge. This means that most of the entries of the adjacency matrix are 0. Therefore, that matrix does not provide a very efficient way of encoding the graph. A more efficient way is provided by simply listing for each  $i$  those vertices that send links to  $i$ , together with the corresponding weights in the weighted case.

An isomorphism between graphs  $\Gamma_1 = (V_1, E_1), \Gamma_2 = (V_2, E_2)$  is a bijection  $\Phi : V_1 \rightarrow V_2$  that preserves neighborhood relations, that is,  $i \sim j$  iff  $\Phi(i) \sim \Phi(j)$ . In other words,  $i$  and  $j$  are connected by an edge precisely if their images under  $\Phi$  are. Isomorphisms preserve the degrees of vertices, that is,  $n_i = n_{\Phi(i)}$  for every vertex  $i$ . An automorphism of  $\Gamma$  is an isomorphism from  $\Gamma$  onto itself. The identity map of the vertex set of  $\Gamma$  is obviously an automorphism, but there may or may not be others, depending on the structure of  $\Gamma$ . The automorphisms of  $\Gamma$  form a group under composition. We can then quantify the symmetry of  $\Gamma$  as the order of its automorphism group.

The number of graphs of order  $k$  grows very fast as a function of  $k$ , and therefore, it becomes unwieldy already for rather small  $k$  to list all graphs of order  $k$ . Therefore, it is of interest to develop constructions for particular classes or types of graphs. There exist deterministic and stochastic construction schemes. We shall discuss stochastic constructions below in 3.5 in the chapter on stochastic processes. Deterministic constructions typically produce rather regular graphs, that is ones with high degrees of symmetries whereas the stochastic constructions can produce typical representatives of larger classes of graphs. A paradigm of a symmetric graph is a complete graph, meaning that any two different vertices are connected by an edge. For a complete graph, every bijection of its vertices yields an automorphism, and therefore, it is maximally symmetric.

A cycle in  $\Gamma$  is a closed path  $i_0, i_1, \dots, i_m = i_0$  for which all the vertices  $i_1, \dots, i_m$  are distinct. For  $m = 3$ , we speak of a triangle. A cycle that contains all the vertices of  $\Gamma$  is called a Hamiltonian cycle (and such a cycle need not exist for a given graph). A graph without cycles is called a tree.

A graph is called  $k$ -regular if all vertices have the same degree  $k$ . As already mentioned, a graph is bipartite if its vertex set can be decomposed into two disjoint components  $V_1, V_2$  such that whenever  $i \sim j$ , then  $i$  and  $j$  are in different components. It is not hard to see that a graph is bipartite iff it does not contain cycles of odd length. In particular, it cannot contain any triangles.

Another useful concept for analyzing graphs is the  $k$ -core. For  $k \in \mathbb{N}$ , the  $k$ -core of a graph  $\Gamma$  is the not necessarily connected maximal subgraph  $H$  of  $\Gamma$  with the property that every vertex of  $H$  has at least  $k$  neighbors in  $H$ , that is, its degree in  $H$  is at least  $k$ . When we exclude the trivial case of an isolated vertex, then  $\Gamma$  itself coincides with its 1-core. When  $\Gamma$  is a tree, already its 2-core is empty. Every cycle of  $\Gamma$  is contained in its 2-core. The core decomposition of  $\Gamma$ , that is, the successive determination of its  $k$ -cores for increasing  $k$ , is a computationally simple way of decomposing the graph.

There exist other parameters that describe certain – more or less – important qualitative properties of graphs. One set of such parameters arises from the metric on the graph generated by the above assignment of length 1 to every edge. The diameter of the graph is the maximal distance between any two of its nodes. As an example how such a parameter can distinguish between typical and non-typical, special graphs, we record that there exists a constant  $c$  with the property that the fraction of all graphs with  $N$  nodes having diameter exceeding  $c \log N$  tends to 0 for  $N \rightarrow \infty$ . Informally expressed, most graphs of  $N$  nodes have a diameter of order  $\log N$ . Thus, graphs with large diameters, like a chain  $i_1 \sim i_2 \sim \dots \sim i_N$  with no other edges, are rare. In the other direction, that is, considering graphs with very small diameters, of course, a fully connected graph has diameter 1. However, one can realize a small diameter already with much fewer edges; namely, one selects one central node to which every other node is connected. In that manner, one obtains a graph of  $N$  nodes with  $N - 1$  edges and diameter 2. Of course, the central node then has a very large degree, namely  $N - 1$ . It is a big hub. Similarly, one can construct graphs with a few hubs, so that none of them has to be quite that big, efficiently distributed so that the diameter is still rather small. Such graphs can be realized as so-called scale free graphs to be discussed below. Another useful quantity is the average distance between nodes in the graph. The property of having a small diameter or average distance has been called the small-world effect.

A rather different quantity is the clustering coefficient that measures how many connections there exist between the neighbors of a node. Formally, it is defined as

$$C := \frac{3 \times \text{number of triangles}}{\text{number of connected triples of nodes}}. \quad (2.1.1)$$

The normalization is that  $C$  becomes one for a fully connected graph. It vanishes for trees and other bipartite graphs.

A triangle is a cycle of length 3. One may then also count the number of cycles of length  $k$ , for integers  $> 3$ . A different generalization consists in considering complete subgraphs of order  $k$ . Here, the complete  $k$ -graph is the graph with  $k$  vertices and links between all  $i \neq j$ . A  $k$ -clique in a graph  $\Gamma$  is a subgraph that is a complete  $k$ -graph. For example, for  $k = 4$ , we would have a subset of 4 nodes that are all mutually connected. One may then associate a simplicial complex to our graph by assigning a  $k$ -simplex to every such complete subgraph, with obvious incidence relations. For example, two such  $k$ -simplices share a  $(k - 1)$ -dimensional face and are called adjacent when the two corresponding complete  $k$ -subgraphs have a complete  $(k - 1)$ -graph in common. This is the basis of topological combinatorics, enabling one to apply tools from simplicial topology to graph theory.

A basic question in the analysis of graphs is the cluster decomposition. That means that we try to find subgraphs, the clusters, that are densely connected inside, but only sparsely connected to the rest of the graph. For example, one can try to disconnect the graph by cutting as few edges as possible, to obtain two large (super)clusters, and then perhaps iterate the process inside

these superclusters to find a finer decomposition. Conversely, one can try to build up the clusters from inside, for example by identifying maximal sets of adjacent  $k$ -cliques, or, equivalently, in the simplicial complex defined above, finding maximal sets of  $k$ -simplices that are connected by  $(k - 1)$ -dimensional faces. Here, the clusters found are typically not disjoint, in contrast to the ones obtained by the edge-cutting methods. Of course, one may then analyze the overlap between those clusters.

Concerning the number of edges needed to disconnect a graph, some insight is provided by the following result of Menger:

**Lemma 2.1.1.** *Let  $V_1$  and  $V_2$  be disjoint subsets of the vertex set of a graph  $\Gamma = (V, E)$ . The minimal number of edges that need to be deleted from  $\Gamma$  in order to disconnect it in such a manner that  $V_1$  and  $V_2$  are in different components is equal to the maximum number of edge-disjoint paths (that is no two paths are allowed to have an edge in common, even though they may well pass through the same vertex) with one endpoint in  $V_1$  and the other in  $V_2$ .*

Another general question is to identify the most important “core” of the graph. The  $k$ -core defined above is one useful concept for that. The idea there is that a node is important when it is connected with other important nodes. Thus, one finds the core by successively deleting the less important nodes. That procedure might make some nodes that have originally been highly connected, that is, have a large degree, less relevant, because they had only been connected to other nodes of low degrees. Therefore, in particular, the degree of a node in general is not a good measure of its importance. One can also quantify the importance of a vertex or an edge by counting how many shortest connections between pairs of nodes pass through them. Again, one should be a bit cautious here because in some cases, there exist alternatives to shortest paths that are not substantially longer but that avoid the vertex or edge in question. In other words, sometimes vertices or edges can easily be replaced as parts of short connections while in other cases that may not be possible. When one decides the importance according to such considerations, this effect should also be taken into account.

### 2.1.3 The graph Laplacian and its spectrum

As before,  $\Gamma$  is a finite and connected graph. Probably the most powerful and comprehensive set of invariants comes from the spectrum of the graph Laplacian of  $\Gamma$  to which we now turn. (In general terms, this means that, in order to analyze a graph  $\Gamma$ , we shall study functions defined on  $\Gamma$ . These functions will then be decomposed in terms of a particular set of basis functions, as in Fourier analysis. From those basis functions, we shall obtain spectral values that incorporate the characteristic properties of  $\Gamma$ .)

There are several non-equivalent definitions of the graph Laplacian employed in the literature. In order to clarify this issue, we assign weights  $b_i (> 0)$  to the

vertices<sup>1</sup> and introduce an  $L^2$ -product for (complex-valued) functions on  $\Gamma$ :

$$(u, v) := \sum_{i \in V} b_i u(i) \bar{v}(i). \quad (2.1.2)$$

(Since we shall only consider real operators below, it suffices to consider real valued functions, and then the complex conjugate in (2.1.2) is not relevant.)

The most natural choices are  $b_i = 1$  or  $b_i = n_i$  where  $n_i$  is the degree of the vertex  $i$ .<sup>2</sup> We may then choose an orthonormal base of that space  $L^2(\Gamma)$ . In order to find such a basis that is also well adapted to dynamical aspects, we study the graph Laplacian

$$\begin{aligned} \Delta : L^2(\Gamma) &\rightarrow L^2(\Gamma) \\ \Delta v(i) &:= \frac{1}{b_i} \left( \sum_{j, j \sim i} v(j) - n_i v(i) \right) \end{aligned} \quad (2.1.3)$$

where  $j \sim i$  means that  $j$  is a neighbor of  $i$ .<sup>3</sup>

We, in contrast to much of the literature on graph theory (see e.g. [14]), but in accordance with [5], prefer the weights  $b_i = n_i$  over  $b_i = 1$  because the former are well adapted to random walks and conservation laws. (When we have a particle randomly moving on a graph with step size 1 then when it is at vertex  $i$  it can choose each of the neighbors of  $i$  with probability  $1/n_i$  for its next move, and this leads to the corresponding factor in the Laplace operator underlying that random walk.)

The idea behind the definition of  $\Delta$  is of course that one compares the value of a function  $v$  at a vertex  $i$  with the average of the values at the neighbors of  $i$ . When that average is larger than the value at  $i$ , we have  $(\Delta v)(i) > 0$ .

The important properties of  $\Delta$  are the following ones:

1.  $\Delta$  is selfadjoint w.r.t.  $(\cdot, \cdot)$ :

$$(u, \Delta v) = (\Delta u, v) \quad (2.1.4)$$

for all  $u, v \in L^2(\Gamma)$ .<sup>4</sup> This holds because the neighborhood relation is symmetric.

<sup>1</sup>These vertex weights should not be confused with the edge weights discussed above; in other words, here, we are *not* considering weighted graphs in the sense defined above.

<sup>2</sup>For purposes of normalization, one might wish to put an additional factor  $N$  in front of the product where  $N$  is the number of elements of the graph or, equivalently, divide all the vertex weights by  $N$ , but we have decided to omit that factor in our conventions.

<sup>3</sup>There are several different definitions of the graph Laplacian in the literature. Some of them are equivalent to ours inasmuch as they yield the same spectrum, but others are not. The reason is simply that the weights in the underlying product are chosen differently. The operator  $\mathcal{L}v(i) := n_i v(i) - \sum_{j, j \sim i} v(j)$  that is often employed in the literature corresponds to the weights  $b_i = 1$  (up to the minus sign, of course). The operator  $\mathcal{L}v(i) := v(i) - \sum_{j, j \sim i} \frac{1}{\sqrt{n_i} \sqrt{n_j}} v(j)$  employed in the monograph [5], apart from the minus sign, has the same eigenvalues as  $\Delta$  for the weights  $b_i = n_i$ : if  $\Delta v(i) = \mu v(i)$ , then  $w(i) = \sqrt{n_i} v(i)$  satisfies  $\mathcal{L}w(i) = -\mu w(i)$ .

<sup>4</sup>An operator  $A = (A_{ij})$  is symmetric w.r.t. a product  $\langle v, w \rangle := \sum_i b_i v(i) \bar{w}(i)$ , that is,  $\langle Av, w \rangle = \langle v, Aw \rangle$  if  $b_i A_{ij} = b_j A_{ji}$  for all indices  $i, j$ . The  $b_i$  are often called multipliers in the literature.

2.  $\Delta$  is nonpositive:

$$(\Delta u, u) \leq 0 \quad (2.1.5)$$

for all  $u$ . This follows from the Cauchy-Schwarz inequality.

3.  $\Delta u = 0$  precisely when  $u$  is constant. This one sees by observing that, when  $\Delta u = 0$ , there can neither be a vertex  $i$  with  $u(i) \geq u(j)$  for all  $j \sim i$  with strict inequality for at least one such  $j$ , that is, a nontrivial local maximum, nor a nontrivial local minimum, as this would contradict the fact that  $\Delta u(i) = 0$  means that the value  $u(i)$  is the average of the values at the neighbors of  $i$ . Since  $\Gamma$  is connected,  $u$  then has to be a constant (when  $\Gamma$  is not connected, a solution of  $\Delta u = 0$  is constant on every connected component of  $\Gamma$ .)

The preceding properties have consequences for the eigenvalues of  $\Delta$ :

- By 1, the eigenvalues are real.
- By 2, they are nonpositive. We write them as  $-\lambda_k$  so that the eigenvalue equation becomes

$$\Delta u_k + \lambda_k u_k = 0. \quad (2.1.6)$$

- By 3, the smallest eigenvalue is  $\lambda_0 = 0$ . Since we assume that  $\Gamma$  is connected, this eigenvalue is simple, that is

$$\lambda_k > 0 \quad (2.1.7)$$

for  $k > 0$  where we order the eigenvalues as

$$\lambda_0 = 0 < \lambda_1 \leq \dots \leq \lambda_K$$

where we put  $K := N - 1$ .

We next consider, for neighbors  $i, j$ ,

$$Du(i, j) := u(i) - u(j). \quad (2.1.8)$$

$D$  can be considered as a map from functions on the vertices of  $\Gamma$  to functions on the edges of  $\Gamma$ . In order to make the latter space also an  $L^2$ -space, we introduce the product

$$(Du, Dv) := \sum_{e=(i,j)} (u(i) - u(j))(v(i) - v(j)). \quad (2.1.9)$$

Note that we are summing here over edges, and not over vertices. If we did the latter, we would need to put in a factor  $1/2$  because each edge would then be counted twice. We also point out that in contrast to the product of (2.1.2),  $(u, v) = \sum_i b_i u(i)v(i)$ , we do not include weights here. The reason is that here the sum should be considered as a sum of edges and not one over vertices, and since we are considering unweighted graphs at this point, the edges do not carry any natural weights.

The product (2.1.9) encodes more information about the graph than the product (2.1.2). The latter only depends on the weights, but not on the connection structure of the graph. There exist many structurally quite diverse graphs with the same weight sequence, and given a graph, one can rewire it by a cross exchange of edges without changing the degrees of the nodes. Namely, given vertices  $i_1 \sim j_1$  and  $i_2 \sim j_2$ , but without edges between  $i_1$  and  $i_2$ , nor between  $j_1$  and  $j_2$ , we create a new graph by deleting the edges between  $i_1$  and  $j_1$  and between  $i_2$  and  $j_2$  and inserting new edges between  $i_1$  and  $i_2$  and between  $j_1$  and  $j_2$ . That operation preserves the degrees of all vertices, and therefore also the product (2.1.2) for any functions  $u, v$  on the graph. (2.1.9), in contrast, is affected because the edge set is changed.

We have

$$\begin{aligned} (Du, Dv) &= \sum_i \frac{1}{2} (n_i u(i)v(i) + \sum_j n_j u(j)v(j) - 2 \sum_{j \sim i} u(i)v(j)) \\ &= - \sum_i u(i) \sum_{j \sim i} (v(j) - v(i)) \\ &= -(u, \Delta v). \end{aligned} \tag{2.1.10}$$

Thus, our product (2.1.9) is naturally related to the Laplacian  $\Delta$ .

We may find an orthonormal basis of  $L^2(\Gamma)$  consisting of eigenfunctions of  $\Delta$ ,

$$u_k, \quad k = 0, \dots, K$$

( $K = N - 1$ ). This is achieved as follows. We iteratively define, with  $H_0 := H := L^2(\Gamma)$  being the Hilbert space of all real-valued functions on  $\Gamma$  with the scalar product  $(\cdot, \cdot)$ ,

$$H_k := \{v \in H : (v, u_i) = 0 \text{ for } i \leq k - 1\}, \tag{2.1.11}$$

starting with a constant function  $u_0$  as the eigenfunction for the eigenvalue  $\lambda_0 = 0$ . Also

$$\lambda_k := \inf_{u \in H_k - \{0\}} \frac{(Du, Du)}{(u, u)}, \tag{2.1.12}$$

that is, we claim that the eigenvalues can be obtained as those infima. First of all, since  $H_k \subset H_{k-1}$ , we have

$$\lambda_k \geq \lambda_{k-1}. \tag{2.1.13}$$

Secondly, since the expression in (2.1.12) remains unchanged when a function  $u$  is multiplied by a nonzero constant, it suffices to consider those functions that satisfy the normalization

$$(u, u) = 1 \tag{2.1.14}$$

whenever convenient.

We may find a function  $u_k$  that realizes the infimum in (2.1.12), that is

$$\lambda_k = \frac{(Du_k, Du_k)}{(u_k, u_k)}. \tag{2.1.15}$$



Since then for every  $\varphi \in H_k, t \in \mathbb{R}$

$$\frac{(D(u_k + t\varphi), D(u_k + t\varphi))}{(u_k + t\varphi, u_k + t\varphi)} \geq \lambda_k, \quad (2.1.16)$$

the derivative of that expression w.r.t.  $t$  vanishes at  $t = 0$ , and we obtain, using (2.1.10)

$$0 = (Du_k, D\varphi) - \lambda_k(u_k, \varphi) = -(\Delta u_k, \varphi) - \lambda_k(u_k, \varphi) \quad (2.1.17)$$

for all  $\varphi \in H_k$ ; in fact, this even holds for all  $\varphi \in H$ , and not only for those in the subspace  $H_k$ , since for  $i \leq k - 1$

$$(u_k, u_i) = 0 \quad (2.1.18)$$

and

$$(Du_k, Du_i) = (Du_i, Du_k) = -(\Delta u_i, u_k) = \lambda_i(u_i, u_k) = 0 \quad (2.1.19)$$

since  $u_k \in H_k$ . Thus, if we also recall (2.1.10),

$$(\Delta u_k, \varphi) + \lambda_k(u_k, \varphi) = 0 \quad (2.1.20)$$

for all  $\varphi \in H$  whence

$$\Delta u_k + \lambda_k u_k = 0. \quad (2.1.21)$$

Since, as noted in (2.1.14), we may require

$$(u_k, u_k) = 1 \quad (2.1.22)$$

for  $k = 0, 1, \dots, K$  and since the  $u_k$  are mutually orthogonal by construction, we have constructed an orthonormal basis of  $H$  consisting of eigenfunctions of  $\Delta$ . Thus we may expand any function  $f$  on  $\Gamma$  as

$$f(i) = \sum_k (f, u_k) u_k(i). \quad (2.1.23)$$

We then also have

$$(f, f) = \sum_k (f, u_k)^2 \quad (2.1.24)$$

since the  $u_k$  satisfy

$$(u_j, u_k) = \delta_{jk}, \quad (2.1.25)$$

the condition for being an orthonormal basis. Finally, using (2.1.24) and (2.1.10), we obtain

$$(Df, Df) = \sum_k \lambda_k (f, u_k)^2. \quad (2.1.26)$$

We next state **Courant's minimax principle**:

Let  $P^k$  be the collection of all  $k$ -dimensional linear subspaces of  $H$ . We have

$$\lambda_k = \max_{L \in P^k} \min \left\{ \frac{(Du, Du)}{(u, u)} : u \neq 0, (u, v) = 0 \text{ for all } v \in L \right\} \quad (2.1.27)$$

and dually

$$\lambda_k = \min_{L \in \mathcal{P}^{k+1}} \max \left\{ \frac{(Du, Du)}{(u, u)} : u \in L \setminus \{0\} \right\}. \quad (2.1.28)$$

In words: In (2.1.27), we consider the minimal Rayleigh quotient under  $k$  constraints, and we maximize that w.r.t. the constraints. In (2.1.28), we consider the maximal Rayleigh quotient for  $k + 1$  degrees of freedom, and we minimize that w.r.t. those degrees of freedom.

To verify these relations, we recall (2.1.12)

$$\lambda_k = \min \left\{ \frac{(Du, Du)}{(u, u)} : u \neq 0, (u, u_j) = 0 \text{ for } j = 0, \dots, k-1 \right\}. \quad (2.1.29)$$

Dually, we have

$$\lambda_k = \max \left\{ \frac{(Du, Du)}{(u, u)} : u \neq 0 \text{ linear combination of } u_j \text{ with } j \leq k \right\}. \quad (2.1.30)$$

The latter maximum is realized when  $u$  is a multiple of the  $k$ th eigenfunction, and so is the minimum in (2.1.29). If now  $L$  is any  $k + 1$ -dimensional subspace, we may find some  $v$  in  $L$  that satisfies the  $k$  conditions

$$(v, u_j) = 0 \text{ for } j = 0, \dots, k-1. \quad (2.1.31)$$

From (2.1.24) and (2.1.26), we then obtain

$$\frac{(Dv, Dv)}{(v, v)} = \frac{\sum_{j \geq k} \lambda_j (v, u_j)^2}{\sum_{j \geq k} (v, u_j)^2} \geq \lambda_k. \quad (2.1.32)$$

This implies

$$\max_{v \in L \setminus \{0\}} \frac{(Dv, Dv)}{(v, v)} \geq \lambda_k. \quad (2.1.33)$$

We then obtain (2.1.28). (2.1.27) follows in a dual manner.

For a fully connected graph, when all the weights  $b_i$  are equal, also all the nontrivial eigenvalues are equal. For our preferred choice of weights,  $b_i = n_i (= N - 1$  for a fully connected graph of  $N$  vertices), we have

$$\lambda_1 = \dots = \lambda_K = \frac{N}{N-1} \quad (2.1.34)$$

since

$$\Delta v = -\frac{N}{N-1}v \quad (2.1.35)$$

for any  $v$  that is orthogonal to the constants, that is

$$\frac{1}{N} \sum_{i \in \Gamma} n_i v(i) = 0. \quad (2.1.36)$$

In more detail, for a fully connected graph of  $N$  vertices, for  $v$  satisfying (2.1.36),

$$\begin{aligned}
\Delta v(i) &= \frac{1}{n_i} \sum_{j, j \sim i} v(j) - v(i) \\
&= \frac{1}{N-1} \sum_{j \neq i} v(j) - v(i) \\
&= \left(-\frac{1}{N-1} - 1\right)v(i) \quad \text{since by (2.1.36) } v(i) = -\sum_{j \neq i} v(j) \\
&= -\frac{N}{N-1}v_i.
\end{aligned}$$

We also recall that since  $\Gamma$  is connected, the trivial eigenvalue  $\lambda_0 = 0$  is simple. If  $\Gamma$  had two components, then the next eigenvalue  $\lambda_1$  would also become 0. A corresponding eigenfunction would be equal to a constant on each component, the two values chosen such (2.1.36) is satisfied; in particular, one of the two would be positive, the other one negative. We therefore expect that for graphs with a pronounced community structure, that is, for ones that can be broken up into two large components by deleting only few edges as discussed above, the eigenvalue  $\lambda_1$  should be close to 0. Formally, this is easily seen from the variational characterization

$$\lambda_1 = \min \left\{ \frac{\sum_{i,j; j \sim i} (v(i) - v(j))^2}{\sum_i b_i v(i)^2} : \sum_i b_i v(i) = 0 \right\} \quad (2.1.37)$$

(see (2.1.12) and observe that  $\sum_i b_i v(i) = 0$  is equivalent to  $(v, u_0) = 0$  as the eigenfunction  $u_0$  is constant). Namely, if two large components of  $\Gamma$  are only connected by few edges, then one can make  $v$  constant on either side, with opposite signs so as to respect the normalization (2.1.36) with only a small contribution from the numerator.

More generally, when  $\Gamma$  consists of several clusters with only very few connections between them, one should find several eigenvalues close to 0.

The strategy for obtaining an eigenfunction for the first eigenvalue  $\lambda_1$  is, according to (2.1.37), to do the same as one's neighbors. Because of the constraint  $\sum_i b_i v(i) = 0$ , this is not globally possible, however. The first eigenfunction thus exhibits oscillations with the lowest possible frequency. Thus, if we take such a first eigenfunction  $u_1$  and consider the connected components that remain after deleting all edges at whose endpoints  $u_1$  has different signs, then there are precisely two such components, one on which  $u_1$  is positive and one on which it is negative. More generally, the number of connected components of  $\Gamma$  where an eigenfunction for the  $k$ th eigenvalue has a fixed sign is at most  $k + 1$  when the eigenvalues are ordered in increasing order and appropriately when they are not simple, according to a version of Courant's nodal domain theorem proved by Gladwell-Davies-Leydold-Stadler [13].

By way of contrast, according to (2.1.28), the highest eigenvalue is given by

$$\lambda_K = \max_{u \neq 0} \frac{(Du, Du)}{(u, u)}. \quad (2.1.38)$$

Thus, the strategy for obtaining an eigenfunction for the highest eigenvalue is to do the opposite what one's neighbors are doing, for example to assume the value 1 when the neighbors have the value -1. Thus, the corresponding eigenfunction will exhibit oscillations with the highest possible frequency. Here, the obstacle can be local. Namely, any triangle, that is, a triple of three mutually connected nodes, presents such an obstacle. More generally, any cycle of odd length makes an alternation of the values 1 and -1 impossible. The optimal situation here is represented by a bipartite graph, that is, a graph that consists of two sets  $\Gamma_+, \Gamma_-$  of nodes without any links between nodes in the same such subset. Thus, one can put  $u_K = \pm 1$  on  $\Gamma_{\pm}$ . For our choice  $b_i = n_i$ , which we shall now adopt for the subsequent discussion, one then finds

$$\lambda_K = 2 \tag{2.1.39}$$

for a bipartite graph.

In contrast, the highest eigenvalue  $\lambda_K$  becomes smallest on a fully connected graph, namely

$$\lambda_K = \frac{N}{N-1} \tag{2.1.40}$$

according to (2.1.36). For graphs that are neither bipartite nor fully connected, this eigenvalue lies strictly between those two extremal possibilities.

Perhaps the following caricature can summarize the preceding: For minimizing  $\lambda_1$  – the minimal value being 0 – one needs two subsets that can internally be arbitrarily connected, but that do not admit any connection between each other. For maximizing  $\lambda_K$  – the maximal value being 2 – one needs two subsets without any internal connections, but allowing arbitrary connections between them. In either situation, the worst case – that is the one of a maximal value for  $\lambda_1$  and a minimal value for  $\lambda_K$  – is represented by a fully connected graph. In fact, in that case,  $\lambda_1$  and  $\lambda_K$  coincide.

Let us consider bipartite graphs in some more detail. We already noted above that on a bipartite graph, we can determine the highest eigenfunction  $u_K$  explicitly, as  $\pm 1$ , being +1 on one set, -1 on the other set of vertices defining the bipartition. In fact, it is clear from that construction that this property is equivalent to the bipartiteness of the graph. Actually, if the graph is bipartite, then even more is true: Whenever  $\lambda_k$  is an eigenvalue, then so is  $2 - \lambda_k$ . Since 0 is an eigenvalue for any graph, this criterion implies our observation that 2 is an eigenvalue. The general statement is not difficult to see: Let  $G_1, G_2$  be the two vertex sets defining the bipartition. When  $u_k$  is an eigenfunction for the eigenvalue  $\lambda_k$ , then

$$\tilde{u}_k(i) := \begin{cases} u_k(i) & \text{for } i \in G_1 \\ -u_k(i) & \text{for } i \in G_2 \end{cases} \tag{2.1.41}$$

is an eigenfunction with eigenvalue  $2 - \lambda_k$  as is readily verified.

We now return to the issue of decomposing a graph by cutting edges. There exists an important relationship of this issue with the first eigenvalue  $\lambda_1$  which

we shall now describe. This is based on a quantity that is analogous to one introduced by Cheeger in Riemannian geometry, but had already been considered earlier in graph theory by Polya. We therefore call it the Polya-Cheeger constant. Letting  $|E|$  denote the number of edges contained in an edge set  $E$ , the Polya-Cheeger constant is

$$h(\Gamma) := \inf \left\{ \frac{|E_0|}{\min(\sum_{i \in V_1} b_i, \sum_{i \in V_2} b_i)} \right\} \quad (2.1.42)$$

where removing  $E_0$  disconnects  $\Gamma$  into the components  $V_1, V_2$ . Thus, we try to break the graph up into two large components by removing only few edges. We may then repeat the process within those components to break them further up until we are no longer able to realize a small value of  $h$ .

We now derive elementary estimates for  $\lambda_1$  from above and below in terms of the constant  $h(\Gamma)$ . Our reference here is [5] (that monograph also contains many other spectral estimates for graphs, as well as the original references). We start with the estimate from above and use the variational characterization (2.1.37). Let the edge set  $E_0$  divide the graph into the two disjoint sets  $V_1, V_2$  of nodes, and let  $V_1$  be the one with the smaller vertex sum  $\sum n_i$ . We consider a function  $v$  that is  $=1$  on all the nodes in  $V_1$  and  $=-\alpha$  for some positive  $\alpha$  on  $V_2$ .  $\alpha$  is chosen so that the normalization  $\sum_{\Gamma} b_i v(i) = 0$  holds, that is,  $\sum_{i \in V_1} b_i - \sum_{i \in V_2} b_i \alpha = 0$ . Since  $V_2$  is the subset with the larger  $\sum b_i$ , we have  $\alpha \leq 1$ . Thus, for our choice of  $v$ , the quotient in (2.1.37) becomes  $\leq \frac{(1+\alpha)^2 |E_0|}{\sum_{i \in V_1} b_i + \sum_{i \in V_2} b_i \alpha^2} = \frac{(\alpha+1)|E_0|}{\sum_{V_1} b_i} \leq 2 \frac{|E_0|}{\sum_{V_1} b_i}$ . Since this holds for all such splittings of our graph  $\Gamma$ , we obtain from (2.1.42) and (2.1.37)

$$\lambda_1 \leq 2h(\Gamma). \quad (2.1.43)$$

The estimate from below is slightly more subtle, and the estimate presented here works only for the choice

$$b_i = n_i. \quad (2.1.44)$$

We consider the first eigenfunction  $u_1$ . Like all functions on our graph, we consider it to be defined on the nodes. We then interpolate it linearly on the edges of  $\Gamma$ . Since  $u_1$  is orthogonal to the constants (recall  $\sum_i n_i u(i) = 0$ ), it has to change sign, and the zero set of our extension then divides  $\Gamma$  into two parts  $\Gamma'$  and  $\Gamma''$ . W.l.o.g.,  $\Gamma'$  is the part with fewer nodes. The points where (the extension of)  $u_1 = 0$  are called boundary points. We now consider any function  $\varphi$  that is linear on the edges, 0 on the boundary, and positive elsewhere on the nodes and edges of  $\Gamma'$ . We also put  $h'(\Gamma') := \inf \left\{ \frac{|E_1|}{\sum_{i \in \Omega} n_i} \right\}$  where removing the edges in  $E_1$  cuts out a subset  $\Omega$  that is disjoint from the boundary. We then

have

$$\begin{aligned}
\sum_{i \sim j} |\varphi(i) - \varphi(j)| &= \int_{\sigma} \#_e(\varphi = \sigma) d\sigma \\
&= \int_{\sigma} \frac{\#_e(\varphi = \sigma)}{\sum_{i: \varphi(i) \geq \sigma} n_i} \sum_{i: \varphi(i) \geq \sigma} n_i d\sigma \\
&\geq \inf_{\sigma} \frac{\#_e(\varphi = \sigma)}{\sum_{i: \varphi(i) \geq \sigma} n_i} \int_{\sigma} \sum_{i: \varphi(i) \geq s} n_i ds \\
&= \inf_{\sigma} \frac{\#_e(\varphi = \sigma)}{\sum_{i: \varphi(i) \geq \sigma} n_i} \sum_i n_i |\varphi(i)| \\
&\geq h'(\Gamma') \sum_i n_i |\varphi(i)|
\end{aligned}$$

when the sets  $\varphi = \sigma$  and  $\varphi \geq \sigma$  satisfy the conditions in the definition of  $h'(\Gamma)$ ; that is, the infimum has to be taken over those  $\sigma < \max \varphi$ . Here,  $\#_e(\varphi = \sigma)$  denotes the number of edges on which  $\varphi$  attains the value  $\sigma$ . Applying this to  $\varphi = v^2$  for some function  $v$  on  $\Gamma'$  that vanishes on the boundary, we obtain

$$\begin{aligned}
h(\Gamma') \sum_i n_i |v(i)|^2 &\leq \sum_{i \sim j} |v(i)^2 - v(j)^2| \\
&\leq \sum_{i \sim j} (|v(i)| + |v(j)|) |v(i) - v(j)| \\
&\leq 2 \left( \sum_i n_i |v(i)|^2 \right)^{1/2} \left( \sum_{i \sim j} |v(i) - v(j)|^2 \right)^{1/2}
\end{aligned}$$

from which

$$\frac{1}{4} h(\Gamma')^2 \sum_i n_i |v(i)|^2 \leq \sum_{i \sim j} |v(i) - v(j)|^2. \quad (2.1.45)$$

We now apply this to  $v = u_1$ , the first eigenfunction of our graph  $\Gamma$ . We have  $h'(\Gamma') \geq h(\Gamma)$ , since  $\Gamma'$  is the component with fewer nodes. We also have

$$\lambda_1 \sum_{i \in \Gamma'} n_i u_1(i)^2 = \frac{1}{2} \sum_{i \in \Gamma'} \sum_{j \sim i} (u_1(i) - u_1(j))^2, \quad (2.1.46)$$

cf. (2.1.15) (this relation holds on both  $\Gamma'$  and  $\Gamma''$  because  $u_1$  vanishes on their common boundary)<sup>6</sup>. (2.1.45) and (2.1.46) yield the desired estimate (under the assumption (2.1.44))

$$\lambda_1 \geq \frac{1}{2} h(\Gamma)^2. \quad (2.1.47)$$

<sup>5</sup>We obtain the factor 1/2 because we are now summing over vertices so that each edge gets counted twice.

<sup>6</sup>To see this, one adds nodes at the points where the edges have been cut, and extends functions by 0 on those nodes. These extended functions then satisfy the analogue of (2.1.10) on either part, as one sees by looking at the derivation of that relation and using the fact that the functions under consideration vanish at those new “boundary” nodes.

From (2.1.43) and (2.1.47), we also observe the inequality

$$h(\Gamma) \leq 4 \tag{2.1.48}$$

for any connected graph, when the weights  $b_i$  are the vertex degrees  $n_i$ .

One can also about the decomposition of a graph by removing vertices instead of edges. This issue is amenable to a similar treatment, and one can define a quantity analogous to  $h(\Gamma)$  that has the number of vertices whose elimination is needed to disconnect the graph in the numerator; see [5] for details.

## 2.2 Descendence relations

### 2.2.1 Trees and phylogenies

Trees are the formal tool for representing ancestor-descendent relations in biology and other fields. At first sight, the concept of a tree as defined below seems not appropriate for that task, however, when one thinks of parent-offspring relationships in sexually recombining species. There, the relationship graph, the so-called pedigree is branching in the backward direction because each individual has two parents, as well as in the forward direction because individuals on average have more than one offspring if the population is not going extinct. When one considers asexual reproduction, however, the situation becomes simpler because each individual then has only one parent, and branching can occur only forward in time when one considers the descendents over the generations of a single ancestor. This, perhaps, is not such an exciting problem, and, in fact, biologists are rather interested in trees for describing phylogenetic relationships between species instead of individuals. The endpoints of a tree, the so-called leaves (see below for the formal definitions), then correspond to a collection of recent species, and one tries to construct a tree in which the internal vertices represent ancestral species that are the common ancestors of all the species below them. Here, one usually assumes that speciation events are binary branchings, that is, one species splits into two daughter species. (In order to make this consistent, at least some biological taxonomists, the cladists, adopt the convention that whenever a new species branches off from an existing one, the remaining part of the latter then is also classified as a new species.) Traditionally, the similarities between species were gauged on the basis of morphological features, and palaeontologists tried to identify the hypothetical ancestral species with ones documented in the fossil record. (In practice, this encounters many problems, but that is not our concern here.) Today, there exists a powerful alternative to that classical method, the comparison on the basis of genetic data. The idea is obvious, to take DNA samples from members of different species and count the differences so as to determine the genetic distances between the species. On the basis of those distances, a hierarchical grouping should be possible that can be represented by a tree. Of course, in practice, this is not so simple. First of all, the genetic samples need to be comparable. For that, one needs to identify DNA segments in the species representatives that are homologous to

each other, that is, derived from the same ancestral sequence through a process of accumulation of mutations. Since besides point mutations in the DNA, there can also occur rearrangements like insertions, deletions, inversions, first the problem of sequence alignment needs to be addressed and solved for the samples at hand. Next, one assumes that mutations occurred at the same rate in the different lineages, the hypothesis of the molecular clock. Otherwise, the number of genetic differences would not be a uniform measure of the time since branching from a common ancestor. Moreover, one needs to find genetic regions that have not been under selective pressure, but rather where there is a uniform probability of the retention of any mutation. Under stabilizing selection, most mutations are eliminated, and this would lead to an underestimate for the time since branching. For directed selection, in contrast, adaptive pressure leads to a more rapid accumulation of mutations and then to an overestimate of the time since branching.

Even if one can align the sequences successfully and eliminate selection effects, there still remain substantial problems. Often, the genetic distances vary with the genomic regions considered. Thus, depending on the DNA region considered, one might get a different tree. In that case, one might try to find some kind of compromise tree. That will depend on the criterion adopted, however, as we shall discuss a little more below. Sometimes, the data even do not fit into a tree because distances on a tree need to satisfy some necessary conditions discussed below. The question then is what substitute to choose for a tree, an issue that we shall also address below. Also, a species is not entirely homogeneous, and there are also genetic differences between the members of the same species (otherwise, evolution could not work by differential selection). Therefore, one needs to gauge intraspecies differences against interspecies ones. Finally, speciation is not an event that takes place at one clearly identifiable point in time, but rather is a gradual process of the accumulation of differences between different populations until reproductive barriers emerge that prevent further genetic mixing between those populations. Here, we need to invoke the species concept of modern biology. A species is defined as a population of organisms that can sexually produce viable and fertile offspring among them. In practice, however, sometimes that relationship is not necessarily transitive. That is, there can exist subpopulations  $A_1, \dots, A_k$  such that individuals from  $A_i$  can reproduce with those of  $A_{i+1}$  for all  $i$ , but the ones from  $A_1$  are no longer able to reproduce with those from  $A_k$ . An example are the races of domestic dogs that range from rather large to very small ones. More generally, for the assembly of phylogenetic trees, species are considered as static ensembles, while in reality speciation is a temporally extended dynamic process inside groups of individuals. (As an aside, some of those population dynamics can be reconstructed on the basis of a statistical analysis of the distribution of alleles in recent populations, in particular from their deviations from equilibria defined by independence hypotheses.) In spite of all these problems, phylogenetic tree reconstructions are a useful tool for many biologists. There is one issue, however, that can really deal a deadly blow to the whole concept of the representation of phylogenies by trees. As L. Margulis emphasized, many genetic changes are not caused by mutations



in inherited genomes, but rather by horizontal gene transfer through viruses and other processes. That, of course, cannot be represented in a tree. On the other hand, over the course of evolution, organisms seem to have developed some protective mechanisms against such horizontal gene insertions, and the relative efficiency of those provides some justification for attempting to represent genetic data in a tree. In the light of all the difficulties mentioned above, it is then necessary to develop methods for finding trees that contain as few as possible hypotheses not supported by the available data.

We now start with the mathematical formalism; we treat a particular class of graphs, the so-called trees. Our basic reference is [33].

A **tree**  $T = (V, E)$  is a graph without cycles.

**Lemma 2.2.1.** *For a graph  $\Gamma = (V, E)$ , the following statements are equivalent:*

1.  $\Gamma$  is a tree, that is, has no cycles.
2. For any two distinct vertices  $i, j$ , there exists a unique path of distinct vertices joining them (we shall call that path a “shortest path” even though we do not yet have specified a metric at this point – it will, however, turn out to be a shortest path for any metric on the tree).
3.  $|V| = |E| + 1$ .
4. The deletion of any edge disconnects  $\Gamma$ .

Since for any graph  $\Gamma = (V, E)$ , we have  $|V| \leq |E| + 1$ , a tree thus is a graph with the minimal number of edges needed to connect a vertex set  $V$ .

The vertices of a tree that have degree 1 are called leaves. The other vertices are called interior vertices. Sometimes, it is convenient to exclude vertices of degree 2. A rooted tree is a tree with one distinguished vertex  $i_0$ , the root.

Rooted trees are the formal tool to represent hierarchical relationships between individual entities. We say that the vertex  $i_1$  is above the vertex  $i_2$ , or in the phylogenetic interpretation to follow that  $i_1$  is an ancestor of  $i_2$ , and  $i_2$  a descendent of  $i_1$ , when the shortest path from  $i_0$  to  $i_2$  passes through  $i_1$ .

In phylogenies, the aim is the comparison between extant species. Those species then are represented as the leaves of some tree, and the rest of the tree then is built with the purpose that the interior vertices represent common ancestors of all the ones below some. Thus, the interior vertices may correspond to hypothetical species on which no data need to be available. Of course, palaeontologists try to identify those interior vertices with fossil species, but the modern data usually consist of genetic data like pieces of DNA sequences for which one rarely has fossil samples. Thus, in palaeontology, it is natural to allow for degree 2 vertices, representing ancestors of a single extant species that are documented in the fossil record. In molecular sequence analysis, however, one would exclude degree 2 vertices because all interior vertices represent hypothetical reconstructions of common ancestors of several descendent species.

In order to proceed with this formalization, we consider  $X$ -trees where  $X$  is some set. In applications,  $X$  of course is a or the data set. An  $X$ -tree is a tree  $T = (V, E)$  together with a map  $\phi : X \rightarrow V$  whose image contains all vertices of degrees 1 and 2. (In the rooted case, we do not require that the root be in the image of  $\phi$  even though it may have degree  $\leq 2$ .) The map need not be injective. For a phylogenetic ( $X$ -) tree, however, we require that  $\phi$  be a bijection onto the leaves of  $T$ . In particular, such a phylogenetic tree has no vertices of degree 2. When every interior vertex has degree 3, we speak of a binary phylogenetic tree. This is a natural assumption in biology because, in evolution, a species can split into two daughter species, and each of those can then split again, and so on, but one does not see the emergence of three or more daughter species at the same time. In fact, much of phylogenetic tree reconstruction is about resolving the question in which temporal order the various splits into daughter species took place.

An  $X$ -split  $A|B$  is a partition of  $X$  into two non-empty subsets  $A, B$ .<sup>7</sup> Thus, in biological applications,  $A$  might represent those members of  $X$  where a certain feature is present, and  $B$  those where that feature is absent. – Two such splits  $A_1|B_1$  and  $A_2|B_2$  are called compatible when at least one of the intersections  $A_1 \cap A_2, A_1 \cap B_2, B_1 \cap A_2, B_1 \cap B_2$  is empty. If, say,  $A_1 \cap B_2 = \emptyset$  then  $A_1 \subset A_2$  and  $B_2 \subset B_1$ , and vice versa, and so, there is an alternative way of expressing compatibility of splits.

When we have an  $X$ -tree  $(T, \phi)$ , then every edge  $e$  of  $T$  induces an  $X$ -split because it decomposes  $T$  into two subgraphs  $T_1, T_2$  (which might include the degenerate case where one of them consists of a single vertex and no edges), and their preimages under  $\phi$  then constitute a split of  $X$ . When we assume that the tree has no vertices of degree 2 – which we shall henceforth do – different edges lead to subgraphs with different leaf sets, and therefore different edges induce different splits of  $X$ . Those splits then are compatible. We denote the splits of  $X$  induced by the  $X$ -tree  $(T, \phi)$  by  $\Sigma(T, \phi)$ , or simply by  $\Sigma(T)$  when the map  $\phi$  is implicitly understood.

The converse question of what classes of splits of  $X$  come from  $X$ -trees is answered by the following result of Buneman

**Theorem 2.2.1.** *Given a collection  $\Sigma$  of  $X$ -splits, there exists an  $X$ -tree  $(T, \phi)$  (which then is unique – up to isomorphism, of course) for which  $\Sigma = \Sigma(T, \phi)$  precisely if all the splits in  $\Sigma$  are pairwise compatible.*

A tree carries an obvious metric, in the sense that we can quantify the distance between vertices  $i_1$  and  $i_2$  by counting the number of edges in the shortest path between them. More generally, we can assign positive weights  $w(e)$  to the edges  $e$  and then take the sum of the weights of the edges in such a path as the distance  $d(i_1, i_2)$ .

When we consider a set  $X$ , there may already exist some distance function on  $X$ , and the question then emerges whether that distance is compatible with the metric on some  $X$ -tree. The answer is pretty simple, and in fact, we can even

---

<sup>7</sup>That  $A$  and  $B$  yield a partition of  $X$  means that  $A \cup B = X$  and  $A \cap B = \emptyset$ .

take something more general than a metric on  $X$ , namely a so-called dissimilarity map, that is, a non-negative map  $\delta : X \times X \rightarrow \mathbb{R}$  with  $\delta(x, x) = 0$  and otherwise positive, and  $\delta(x, y) = \delta(y, x)$  for all  $x, y \in X$ . For example,  $\delta(x, y)$  could just count in how many characters (see below for a formal definition) the elements  $x$  and  $y$  differ.

The question then is whether we can find an  $X$ -tree  $(T, \phi)$  with weights  $w(e)$  on its edges and associated distance function  $d(., .)$  such that

$$\delta(x, y) = d(\phi(x), \phi(y)) \quad (2.2.49)$$

for all  $x, y \in X$ . In that case, we call  $\delta$  a tree metric. The answer is

**Theorem 2.2.2.** *A dissimilarity map  $\delta$  on  $X$  is a tree metric precisely if it satisfies the 4-point condition*

$$\delta(x, y) + \delta(z, w) \leq \max(\delta(x, z) + \delta(y, w), \delta(x, w) + \delta(y, z)) \quad (2.2.50)$$

for all  $x, y, z, w \in X$ .

In the sequel, (2.2.50) will give rise to two different issues. One is whether it holds or not for all points, and this issue is exemplified in the case where  $\delta$  is the metric coming from a quadrilateral graph where  $x, w, y, z$  are arranged in cyclic order, for example  $\delta(x, w) = \delta(w, y) = \delta(y, z) = \delta(z, x) = 1$  and  $\delta(x, y) = \delta(z, w) = 2$ . Thus, (2.2.50) is not satisfied here. The other issue arises when (2.2.50) is satisfied for all quadruples and consists in the question under which conditions we have even strict inequality for certain quadruples.

Since every edge  $e$  of an  $X$ -tree corresponds to a split  $\sigma$  of  $X$ , we can write a tree metric as

$$d = \sum_{\sigma \in \Sigma(T)} w(e_\sigma) \delta_\sigma \quad (2.2.51)$$

where  $e_\sigma$  is the edge inducing the split  $\sigma$  and

$$\delta_\sigma(i, j) = \begin{cases} 1 & \text{if } i, j \text{ are in different components of } T - e \\ 0 & \text{otherwise.} \end{cases}$$

The point here is that the edges  $e_\sigma$  occurring for  $d(x, y)$  in (2.2.51) with  $\delta_\sigma(x, y) = 1$  are precisely those contained in the shortest path from  $x$  to  $y$ .

This will now lead us to the decomposition theorem of Bandelt and Dress[1]. Let  $\delta$  be a dissimilarity map on  $X$ . For a split  $\sigma = A|B$  of  $X$ , we consider

$$i_\delta(\sigma) := \frac{1}{2} \min_{a_1, a_2 \in A, b_1, b_2 \in B} (\max(\delta(a_1, b_1) + \delta(a_2, b_2), \delta(a_1, b_2) + \delta(a_2, b_1)) - (\delta(a_1, a_2) + \delta(b_1, b_2))). \quad (2.2.52)$$

It is not required that the points  $a_1$  and  $a_2$  or  $b_1$  and  $b_2$  be different. For example, this expression can become negative when  $\delta$  does not satisfy the triangle inequality: take  $a_1 = a_2 =: a$  and  $b_1, b_2$  with  $\delta(b_1, b_2) > \delta(b_1, a) + \delta(a, b_2)$ . – In order to understand the significance of  $i_\delta(\sigma)$  better, we consider some examples.

These examples will be graphically displayed in the figure below. We first take a space  $X = \{x, y, z, w\}$  consisting of 4 points, with the condition

$$\delta(x, y) + \delta(z, w) < \max(\delta(x, z) + \delta(y, w), \delta(x, w) + \delta(y, z)). \quad (2.2.53)$$

If  $\delta(x, y) = \delta(z, w) = 2$ ,  $\delta(x, z) = \delta(x, w) = \delta(y, z) = \delta(y, w) = 3$ , the split  $\{x, y\}|\{z, w\}$  has index  $i_\delta = 1/2$  and is induced from a tree with leaves  $x, y, z, w$  and interior nodes  $\xi, \zeta$  with  $\delta(x, \xi) = \delta(y, \xi) = \delta(\xi, \zeta) = \delta(z, \zeta) = \delta(w, \zeta) = 1$ . The splits  $\{x, z\}|\{y, w\}$  or  $\{x, w\}|\{y, z\}$ , however, have  $i_\delta(\sigma) = 0$  and are not induced by that tree metric. When we have equality in (2.2.53), say,  $\delta(x, y) = \delta(z, w) = 2$ ,  $\delta(x, z) = \delta(z, w) = \delta(y, z) = \delta(y, w) = 2$ , the metric can still be represented by a tree metric, this time with a single interior vertex  $\xi$  that has distance 1 from all leaves. We call this tree a star. Here, there is no longer a natural grouping of the vertices into two pairs. – When instead  $\delta(x, y) = \delta(z, w) = \delta(x, z) = \delta(y, w) = 3$ , and  $\delta(x, w) = \delta(y, z) = 4$ , then (2.2.53) holds again. This time, we can represent the metric by a graph with 4 interior vertices  $\xi, \eta, \zeta, \omega$  that is not a tree.  $\xi, \eta, \zeta, \omega$  form a rectangle with  $\delta(\xi, \eta) = \delta(\xi, \zeta) = \delta(\omega, \zeta) = \delta(\eta, \omega) = 1$ , the other nontrivial distances between them being equal to 2, and with  $x$  connected to  $\xi$ ,  $y$  to  $\eta$ ,  $z$  to  $\zeta$ ,  $w$  to  $\omega$ , all with distance 1. Thus, we need to insert an interior rectangle in order to represent the metric on a graph. That rectangle then expresses the ambiguity in the dissimilarity map for a hierarchical grouping. Of course, the rectangle is in fact a square, and so there is some special symmetry. We therefore also consider the case where  $\delta(x, y) = \delta(z, w) = 3$ ,  $\delta(x, z) = \delta(y, w) = 4$ ,  $\delta(x, w) = \delta(y, z) = 5$ . In that case, we again insert 4 interior vertices  $\xi, \eta, \zeta, \omega$  that form a rectangle, this time with  $\delta(\xi, \eta) = \delta(\eta, \omega) = 1$ ,  $\delta(\xi, \zeta) = \delta(\omega, \zeta) = 2$ . In any case, when we have such a rectangle, we produce splits by cutting pairs of parallel edges. Cutting the edges between  $\xi$  and  $\zeta$  and between  $\eta$  and  $\omega$ , for example, produces the split  $\{x, y\}|\{z, w\}$ . Cutting the edges between  $\xi$  and  $\eta$  and between  $\zeta$  and  $\omega$  instead produces the split  $\{x, z\}|\{y, w\}$ . Now, in contrast to the tree case, both these splits have  $i_\delta(\sigma) > 0$ . The split  $\{x, w\}|\{y, z\}$ , however, has  $i_\delta(\sigma) < 0$ . In the tree case, interchanging  $x$  with  $y$  or  $z$  with  $w$  would not have made any difference for the distances between those 4 vertices, but this is no longer so in the rectangle case.

After this example, let us return to the general case. When  $\delta$  is a tree metric from an  $X$ -tree  $(T, \phi)$ , the split  $\sigma$  of  $X$  is induced by that  $X$ -tree precisely if  $i_\delta(\sigma) > 0$ . In that case, we then have  $i_\delta(\sigma) = w(e_\sigma)$  for the weight of the edge inducing the split. And we can rewrite (2.2.51) then as

$$\delta = \sum_{\sigma \text{ } X\text{-split with } i_\delta(\sigma) > 0} i_\delta(\sigma) \delta_\sigma. \quad (2.2.54)$$

The split decomposition theorem of Bandelt and Dress[1] then says that every dissimilarity map can be written as a sum over such tree metrics plus a remainder that has no splits with  $i_\delta(\sigma) > 0$ :

**Theorem 2.2.3.** *Let  $\delta$  be any dissimilarity map on  $X$ . We then have a decom-*

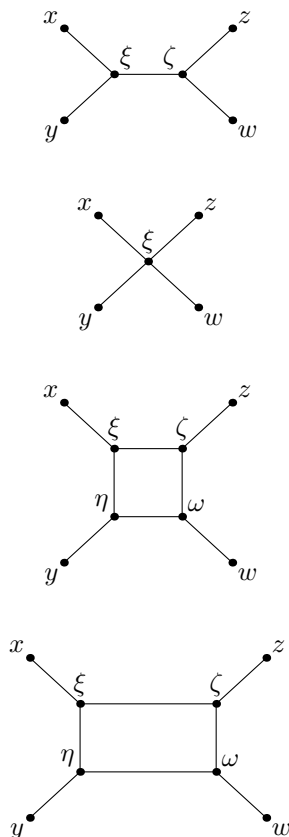


Figure 2.1:

*position*

$$\delta = \delta_0 + \sum_{\sigma \text{ X-split with } i_\delta(\sigma) > 0} i_\delta(\sigma) \delta_\sigma \quad (2.2.55)$$

where  $\delta_0$  admits no splits with  $i_\delta(\sigma) > 0$ .

The star in our above example admits no splits into pairs of points with  $i_\delta(\sigma) > 0$ . This is an undesirable situation in phylogenetic tree reconstruction because the grouping of the four vertices into pairs is ambiguous. However, when we split off a single point from the remaining three, we get  $i_\delta(\sigma) > 0$ . The simplest example of a metric space admitting no splits at all with  $i_\delta(\sigma) > 0$  is given by 5 points  $x, y, z, w, v$  with  $d(x, v) = d(y, z) = d(z, w) = d(y, w) = 2$  and the other distances between different points all being one. To describe this metric space somewhat differently, we take the two sets  $A := \{x, v\}$ ,  $B := \{y, z, w\}$  and connect each point in  $A$  with every point in  $B$  by an edge of length 1. Thus, we see that the graph constructed in this way is bipartite – in the terminology of graph theory, this is the bipartite graph  $K_{2,3}$ . – For this example, then  $\delta_0$  is

nontrivial, and moreover,  $d = \delta_0$ .

When, conversely,  $\delta_0$  vanishes, the dissimilarity map  $\delta$  is called totally decomposable. We recall that in the above example with the interior rectangle, the splits  $\{x, y\}|\{z, w\}$  and  $\{x, z\}|\{y, w\}$  both have positive  $i_\delta(\sigma)$ , and they decompose the metric. Thus, a totally decomposable metric need not be a tree metric. The problem of this example for phylogenetic tree reconstruction is that there is no unique split that decomposes the dissimilarity map, that is, on the basis of the dissimilarity map, we do not know how to group the elements. Another, larger, example of this type, that is, of a totally decomposable metric that is not a tree metric and where therefore the groupings of the elements are not unique, is displayed in the next figure.

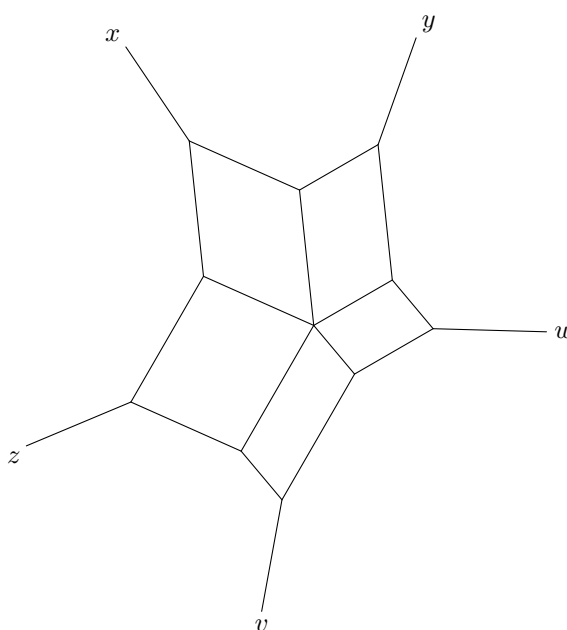


Figure 2.2:

Bandelt and Dress[1] proved that a dissimilarity map  $\delta$  is totally decomposable iff for all  $x, y, z, v, w \in X$ ,

$$i_\delta(\{x, y\}|\{z, v\}) \leq i_\delta(\{x, w\}|\{z, v\}) + i_\delta(\{x, y\}|\{z, w\}). \quad (2.2.56)$$

Splits are decompositions of  $X$  into two subsets. More generally, we can consider characters, that is, functions  $\chi : X_0 \rightarrow S$  where  $\emptyset \neq X_0 \subset X$  and  $S$  is a finite set, the set of character states.  $\chi$  is called non-trivial if there are at least two character states that are each assumed by more than one element

of  $X_0$ . We say that the character  $\chi$  factors through the  $X$ -tree  $(T, \phi)$  when there exists  $\chi' : T \rightarrow S$  (here, we mean by a function on the tree  $T$  a function that is defined on the vertices of  $T$ ) with  $\chi = \chi' \circ \phi|_{X_0}$ . Such a character  $\chi$  that factors through the  $X$ -tree  $(T, \phi)$  is called convex on  $T$  if for each  $a \in S$ , the subgraph with vertex set  $(\chi')^{-1}(a)$  is connected. This is equivalent to the existence, for every pair  $a, b$  of different character states, of an  $X$ -split  $A|B$  of  $T$  with  $(\chi')^{-1}(a) \subset A$  and  $(\chi')^{-1}(b) \subset B$ .

The concept of character convexity is fundamental for the phylogenetic systematics developed by W.Hennig, the so-called cladism. There, one wants to identify monophyletic groups, that is rooted subtrees of phylogenetic trees that contain all the descendants of that vertex that is declared to be the common ancestor and made the root of the subtree. For example, in standard zoological systematics, vertebrates constitute a monophyletic group while fish don't because the other vertebrate groups (amphibians, reptiles, birds, and mammals) are also descendants of fish; in fact, here only birds and mammals are monophyletic in the sense of cladism. We consider an ancestral species  $A$  with daughter species  $A_1$  and  $A_2$ .<sup>8</sup> Of course, this can be represented by a tree with root  $A$  and leaves  $A_1, A_2$ . Suppose that a character state  $a$  in  $A$  is preserved in  $A_1$ , but changed into  $a'$  in  $A_2$ . Now suppose that that species  $A_2$  further splits into two daughter species  $A_{21}$  and  $A_{22}$ . We then get a new tree with leaves  $A_1, A_{21}, A_{22}$  if there are no further splittings while  $A_2$  now is an interior node of degree 3. We consider two cases as displayed in the following figure. In

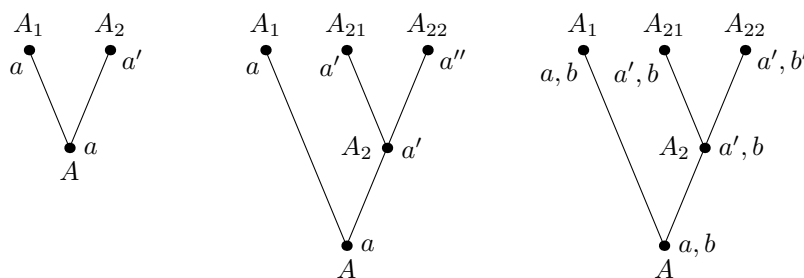


Figure 2.3:

the first case,  $A_{21}$  preserves the state  $a'$  while in  $A_{22}$  it is further transformed into  $a''$ . In the other case, both of them preserve  $a'$ , but in  $A_{22}$  the state of some other character is transformed into the value  $b'$  from the common value  $b$  shared by  $A, A_1, A_2, A_{21}$ . In such a situation, the ancestral states  $a, b$  are called plesiomorph, the derived states  $a', a'', b'$  apomorph. These are relative concepts because  $a'$  is plesiomorph compared with  $a''$ , that is, when we only consider the subtree with root  $A_2$  and leaves  $A_{21}, A_{22}$ . Two species sharing the same

<sup>8</sup>It is a basic principle of cladism that whenever a new species splits off from some line, the remaining part of that line is also classified as a new species. This makes the systematics amenable to tree representations. Moreover, from the morphological approach underlying cladism that is based on paleontological data, any two species differ in the state of at least one character.

plesiomorph state of a character are called symplesiomorphic w.r.t. that character, those sharing an apomorphic state are called synapomorphic. In the last example,  $A_1$  and  $A_{21}$  are symplesiomorphic for  $b$  while  $A_{21}$  and  $A_{22}$  are synapomorphic w.r.t.  $a'$ . In the preceding example, where  $A_{22}$  had the character state  $a''$ , the states  $a', a''$  together constitute a synapomorphy between  $A_{21}$  and  $A_{22}$ . Only synapomorphy, but not symplesiomorphy, can be an indication of a monophyletic group. Here then enters the convexity assumption. Namely, in order to be able to use shared derived characters, that is, synapomorphies for identifying monophyletic groups, we must exclude the following two possibilities:

1. Reversion: In the last example,  $A_{22}$ , instead of assuming the new state  $a''$ , reverts to the ancestral state  $a$ .
2. Convergence: In the same example,  $A_1$ , instead of keeping the state  $b$ , assumes the same state  $b'$  that originated in the species  $A_{22}$  while  $A_{21}$  kept  $b$ .

Of course, there exist biological examples for either possibility. Snakes have lost the limbs that their ancestors had gained. Birds, bats, and insects have independently developed wings. In fact, the wings of birds and bats are plesiomorph when considered as forelimbs, but not as wings. Sometimes, the distinction between plesiomorphy and apomorphy is not clear or needs to be reconsidered in the light of genetic sequence data. For example, it had been thought for a long time that the eyes in arthropods, molluscs, and vertebrates are examples of a convergent evolution. It has been discovered, however, that eye formation in all these lineages is directed by the same type of controlling gene from the class of homeotic (Hox) genes. An uncontroversial<sup>9</sup> example of convergence is mimikry where one species imitates the coloration or other pattern of an unrelated species that is avoided by predators. In any case, reversion and convergence are relatively rare in biological evolution, however. Both these possibilities are excluded by character convexity.

When one has several characters, one wants to find a single  $X$ -tree for which all of them are convex. When such a tree exists, these characters are called compatible. As for compatibility of splits, there exists a theorem characterizing the compatibility of characters, but since the formulation is more complicated we refer to [33].

When working with biological data, typically not all the characters are compatible, and one then wishes to quantify that non-compatibility and construct a tree that comes as close as possible to rendering all the characters convex. This is the idea of parsimony. More precisely, given a function  $f$  on the vertex set  $V$  of a graph  $\Gamma$ , the changing number of  $f$  is the number of those edges of  $\Gamma$  on whose endpoints  $f$  assumes different values, that is, the number of all edges  $e = (i, j)$  with  $f(i) \neq f(j)$ . Let now  $\chi : X_0 \rightarrow S$  be a character that factors

---

<sup>9</sup>at least as long as one does not look at the underlying genetic mechanisms; in fact, it may well turn out in a given example that the imitation of a pattern is produced by the same kind of genetic regulatory mechanism as the imitated pattern, and at least the general framework of that genetic regulation might be derived from some common ancestor



through the  $X$ -tree  $(T, \phi)$ , with  $\chi = \chi' \circ \phi|_{X_0}$  as above. Here, we are assuming that  $\chi'$  is already defined on all the vertices of  $T$ . Of course, it is then arbitrary how to define  $\chi'$  on those vertices of  $T$  that are not in the image of  $\phi(X_0)$ , in case  $\phi$  is not surjective on  $X_0$ . For a character  $\chi$ , we then define its parsimony score  $s(\chi, T)$  for the  $X$ -tree  $(T, \phi)$  as the minimal changing number of all those extensions  $\chi'$  on  $T$  that factor  $\chi$ . Given a set of characters, its parsimony score on an  $X$ -tree then is simply the sum of the individual parsimony scores. A maximal parsimony  $X$ -tree for that set of characters then is one that minimizes that parsimony score.

It is not difficult to see that the parsimony score is related to character convexity. In fact, given a character that assumes  $\nu$  different states, the so-called homeoplasy of the character  $\chi$  on  $T$

$$h(\chi, T) := s(\chi, T) - \nu + 1 \geq 0 \quad (2.2.57)$$

with equality precisely if  $\chi$  is convex on  $T$ . Thus, the total homeoplasy of a character set, the sum of the individual homeoplasies, is also non-negative and vanishes precisely when the characters are compatible.

The concept of maximum parsimony trees is not without difficulties, both conceptually and mathematically. The conceptual difficulties arise from the arbitrariness in the definition and choice of characters. It is a fundamental problem in paleontology and morphology to clearly state what a character is and to decide which characters are independent of each other. Of course, large sets of dependent characters would bias the parsimony concept. The mathematical problems will be taken up below when we consider stochastic processes on trees and other graphs. We shall then see that any method of reconstructing a structure from a data set depends on a model for the underlying process that created the data.

A standard problem is to amalgamate phylogenetic relationships between subsets of  $X$  as expressed in trees into an encompassing tree representing all of  $X$ . Of course, the issue of compatibility will arise again. The smallest meaningful subsets here consist of 4 elements and are called quartets, and trees with 4 leaves are called quartet trees. Also, if one has data about the relationships between the elements of  $X$  and wants to construct a tree or, more generally, find out whether these relationships fit into a tree, a natural strategy is to first construct all local quartet trees and then assemble those into a common tree. When we have a collection  $\mathcal{Q}$  of quartet trees that contains exactly one quartet tree  $\{a, b\}|\{c, d\}$  for every quartet  $Y = \{a, b, c, d\}$  of  $X$ , then, as discovered by Colonius and Schulze [6], there exists a unique  $X$ -tree containing all these quartet tree iff the following two quartet rules hold for all  $a, b, c, d, e \in X$ :

1.

If  $\{a, b\}|\{c, d\}, \{a, b\}|\{d, e\} \in \mathcal{Q}$ , with  $c \neq e$ , then  $\{a, b\}|\{c, e\} \in \mathcal{Q}$

2.

If  $\{a, b\}|\{c, d\}, \{a, c\}|\{d, e\} \in \mathcal{Q}$ , then  $\{a, b\}|\{c, e\} \in \mathcal{Q}$ .

In practice, of course, these rules will be violated for some quintuples of elements of  $X$ , and one therefore cannot construct a tree.

There are other, in fact infinitely many, quartet rules. If  $\mathcal{Q}$  does not contain a quartet tree for every quartet in  $X$ , that is, if we only have a subcollection of quartet trees, then we need to invoke more of those rules to check for compatibility, see [33] for more on this topic. For an algorithm for the (re)construction of a tree from quartets, see [35].

### 2.2.2 Genealogies (pedigrees)

While species can be considered as important biological entities in their own right, the ancestor-descendent relationships in phylogenetic trees can also be viewed as accumulated genealogies of the individuals constituting the populations underlying the species. Thus, let us consider those genealogies a little, even though they in turn can be viewed as combinations of inheritance processes of genes passed on from parents to offspring. The latter, in fact, will lead us back to trees below.

The genealogy or pedigree of an individual in a sexually recombining population is a directed graph. Each individual has two incoming links from its parents while the number of outgoing links counts its offspring. Since no individual can be a descendent of its own offspring, or an ancestor of its own parents, the graph is acyclic (it has no directed cycles; the underlying undirected graph may well have cycles as the result of inbreeding in the population). The nodes without outgoing links represent those individuals that did not produce or have not yet produced offspring. When the graph represents a population history, one can essentialize it by pruning all the vertices without outgoing links that corresponds to individuals no longer alive. This will then be an iterative process because in the next step one would have to prune those vertices that have outgoing links only to vertices that have been pruned in the previous step. In that manner, one iteratively eliminates all vertices that do not have living descendents. Thus, one is left with the ancestral relationships leading to the present population.

Since one does not want to extend the pedigree to the infinite past, one starts with some ancestral population. The essentialized pedigree then contains only those members of the ancestral population that have descendents in the present generation. If one moves further to the next generation, then some of those ancestors may cease to have descendents and therefore will get eliminated. Some of those ancestors, called the lucky ones, however, will turn out to be ancestors of all members of the present population, and they will therefore also leave descendents in all future generations, until the entire population goes extinct.

Often, one assumes that the different generations do not overlap. The generations can then be labelled by their distance from the ancestral one, and links always go from generation  $n$  to generation  $n + 1$ .

Also, from the pedigree graph, one can construct another graph expressing mating relationships. In that graph, there is an (undirected) edge between two individuals when they have produced offspring together. When the species is dioecious, that is, has separate sexes, the mating graph is bipartite, the two

classes corresponding to the females and the males. The mating graph usually is not connected, however, therefore, strictly speaking, violating our definition of a graph. When the population is strictly monogamous, the graph consists of disjoint pairs only, after we have essentialized it and eliminated all the bachelors and spinsters.

Of course, this is all rather simple. Later on, when we consider stochastic branching processes, pedigrees of sexually recombining populations become rather difficult, but for the moment we leave the subject and turn to

### 2.2.3 Gene genealogies (coalescents)

The pedigree just considered for a dioecious population contains two trees (more precisely, so called forests, that is, not necessarily connected unions of trees) as subgraphs, namely the ones corresponding to the male and the female individuals. Let us take one of them, say the female one. Thus, we only consider mother-daughter relationships. For two individuals, we can then ask when their lineages coalesce or merge back in the past, that is, how many generations back they had the same female ancestor. For two sisters, we need only go one generation back, as they have the same mother, while first (in the female line) cousins share a maternal grandmother, and so on. Once the lineages coalesce, they will stay together all the way back to the ancestral population. Of course, in principle, they may never merge, that is the two females under consideration may be descendents of different females in the ancestral population. When we go sufficiently many generations back into the past, however, with overwhelming probability, all presently living females in the populations will descend from the same ancestral female, the “Eve”. All other females in that ancestral population will then have no descendents from an uninterrupted female line in the present populations; of course they may or may not have descendents from some lineages that include some males. As already described above, we can essentialize the graph by eliminating all females without female descendents in an iterative manner so that only those remain that have an uninterrupted line of female descendents down to the present sample. When we do coalescence theory, that is, follow the ancestry of the present sample back in time, then, in fact, those eliminated individuals will never occur in the consideration. This represents an enormous simplification in practice when compared with considering the forward branching process for the (female) descendents of an ancestral populations where all descendents will occur regardless of whether they contribute to future generations or not.

Let us consider this scenario in more detail in a simple example that will lead us to the Wright-Fisher model of population genetics. We consider a population with non-overlapping generations, and we assume that the size of the population remains constant  $= 2N$  across generations. We also assume, for simplicity, that the sex ratio remains constant and equal so that we are dealing with a population of  $N$  females. The assumption of the Wright-Fisher model is that, given generation  $n$ , consisting of a population of  $N$  individuals, generation  $n + 1$

is (mathematically) created by choosing  $N$  times randomly and independently an individual from generation  $n$  as mother. Here, it is assumed that the population is entirely homogeneous, or, in more biological terms, that all members are equally fit, so that at each selection step, each member has the same chance of being chosen. Also, creating daughters does not affect the fitness, and so, the chance to be chosen at a given step does not depend on how often one has already been chosen in previous steps. Putting it another way, each individual in generation  $n + 1$  picks individual  $j$  in generation  $n$  with probability  $1/N$  as its mother, and this sampling is carried out  $N$  times with replacement. If  $d_j$  is the number of daughters of individual  $j$ , we thus have for the probability of having  $\nu$  daughters

$$p(d_j = \nu) = \binom{N}{\nu} \left(\frac{1}{N}\right)^\nu \left(1 - \frac{1}{N}\right)^{N-\nu}. \quad (2.2.58)$$

This is a binomial distribution,  $Bi(N, \frac{1}{N})$ , and so, the number of daughters of a given female is binomially distributed. The expectation value is

$$E(d_j) = N \frac{1}{N} = 1 \quad (2.2.59)$$

which of course reflects the fact that the population size is constant, and the variance is

$$Var(d_j) = N \frac{1}{N} \left(1 - \frac{1}{N}\right) = 1 - \frac{1}{N} \quad (2.2.60)$$

(see (3.1.20) below). The correlation between the numbers of daughters of different females  $j, k$  is

$$Cor(d_j, d_k) = \frac{Cov(d_j, d_k)}{\sqrt{Var(d_j)Var(d_k)}} = -\frac{1}{N-1}. \quad (2.2.61)$$

The correlation is negative, again because the population size is constant, and therefore, when  $j$  has many daughters, there is less room for  $k$  to have many daughters as well (when we already know that an individual different from  $j$  has one daughter, then the expected number of daughters of  $j$  is reduced to  $\frac{N-1}{N}$  in place of the value 1 of (2.2.59)). This effect is rather small in large populations. For large  $N$ , the binomial distribution  $Bi(N, \frac{1}{N})$  is approximated by a Poisson distribution

$$p(d_j = \nu) \approx \frac{1}{\nu!} e^{-1} \quad (2.2.62)$$

with mean and variance = 1 (see (3.1.8), (3.1.9) in 3.1 below). In particular, the probability of having no daughters is

$$p(d_j = 0) \approx e^{-1} \approx .37 \quad (2.2.63)$$

while then the probability to have at least one daughter becomes

$$p(d_j > 0) \approx 1 - e^{-1} \approx .63 \quad (2.2.64)$$

Therefore, the present population descends from a fraction of about  $.63^n$  females  $n$  generations ago. Of course, this eventually goes to 0 for large  $n$  which leads to the absurd result that the present females derive from fewer than one individual in the ancestral generation. Of course, the puzzle is resolved by observing that these approximations were only valid for large population sizes. For small populations, a more refined analysis is needed. This is the subject of coalescence theory, originally founded by J.Kingman.

Again, we stay with our simple example and ask for the distribution of the number  $T_2$  of generations that we need to go back in time to find a common ancestor of two individuals from the present population. That is, we seek the time to the most recent common ancestor (MRCA) of the two individuals. The probability that the two individuals have the same mother, that is, that the MRCA is found already in the first generation from the past, is  $\frac{1}{N}$  because once we have identified the mother of the first individual, the probability that the second one has the same mother is  $\frac{1}{N}$ . Thus, the two have different mothers with probability  $1 - \frac{1}{N}$ . Iteratively, the chance to find the MRCA  $n$  generations back then is

$$p(T_2 = n) = \left(1 - \frac{1}{N}\right)^{n-1} \frac{1}{N} \quad (2.2.65)$$

because they then have different ancestors in  $n - 1$  generations. This is a geometric distribution, and its mean is

$$E(T_2) = \frac{1}{\frac{1}{N}} = N \quad (2.2.66)$$

which is equal to the population size.

In a similar manner, we can consider the time to find the MRCA for  $M$  individuals. The probability that  $m$  individuals have all different mothers is

$$\frac{N-1}{N} \frac{N-2}{N} \cdots \frac{N-m+1}{N} = \prod_{\mu=1}^{m-1} \left(1 - \frac{\mu}{N}\right) = 1 - \binom{m}{2} \frac{1}{N} + O\left(\frac{1}{N^2}\right) \quad (2.2.67)$$

because when the mother of the first individual is determined, there are  $N - 1$  possibilities for the second to have a different mother, and when that is also determined, there remain  $N - 2$  possibilities for the third individual to have a mother different from the previous two, and so on. Thus, neglecting terms of order  $\frac{1}{N^2}$  for a large population size  $N$ , a coalescence event occurs in a given generation with probability  $\binom{m}{2} \frac{1}{N}$ , while no coalescence event occurs with probability  $1 - \binom{m}{2} \frac{1}{N}$ . Thus, the probability distribution for the time  $T_m$  of a coalescence event that reduces the number of different ancestors from  $m$  to  $m - 1$

$$p(T_m = n) = \left(1 - \binom{m}{2} \frac{1}{N}\right)^{n-1} \binom{m}{2} \frac{1}{N}. \quad (2.2.68)$$

In analogy to (2.2.66), we have

$$E(T_m) = \frac{1}{\binom{m}{2} \frac{1}{N}} = \frac{2N}{(m-1)m}. \quad (2.2.69)$$

When we then want to go back from  $M$  individuals to a single ancestor, we need to consider all the coalescent events from  $m$  to  $m - 1$  for  $m = 2, \dots, M$ . Since the times for these events are independent of each other, the expected number of generations back in the past for  $M$  female individuals to have a single female ancestor is

$$\sum_{m=2}^M E(T_m) = \sum_{m=2}^M \frac{1}{\binom{m}{2} \frac{1}{N}} = \sum_{m=2}^M \frac{2N}{(m-1)m} = 2N \left(1 - \frac{1}{M}\right). \quad (2.2.70)$$

In other words, this is the expected height (measured in number of generations) of the tree starting with a single female ancestor and leading to the present ensemble of  $M$  females. When we compare (2.2.70) with (2.2.66), we see that the latter is less than 2 times the former. This means that the final step of reducing the number of ancestors from 2 to 1 typically takes at least half the time of the whole process. Thus, the long branches of the tree arise when there are only few females in the ancestry of the sample.

One can, of course, perform the same analysis with males in place of females. Let us insert a small variation, however, to account for the fact that in many animal species, like most mammals, and also in many human societies, the variance in the number of offspring for males is considerably higher than for females while obviously the expectation value is the same, assuming that the population is in gender equilibrium. This higher variance is easy to achieve in our model. The simplest version just stipulates that in each generation only a certain fraction  $0 < q < 1$  of the number of males is having offspring at all. When we then look for the father of an individual, each of those ones is taken with probability  $1/qN$  while the other ones are simply discarded. Thus, two individuals now stand a chance of  $1/qN$  of having the same father. Thus,  $N$  gets replaced by  $qN$  in all formulae. In particular, the expectation values for the waiting times in (2.2.66), (2.2.70) are shortened by a factor  $q$ , and we expect to find the MRCA in the male line correspondingly fewer generations ago than the one in the female line. In other words, “Adam” lived many generations after “Eve”.

Coalescence theory is mainly interested in describing the ancestry of genes, or more precisely, of DNA segments, rather than of individuals. Formally, the basic scenario is the same, however, and therefore, we have described the basic situation above for the more intuitive case of individuals. The basic scenario neglects the issues of mutation and recombination. In order to exclude recombination, there are two possibilities, one of significance for biological data, these other one solely for modeling purposes. The first one consists in considering those DNA segments that do not recombine. One class is given by non-nuclear mitochondrial DNA that is only contained in egg cells, but not in sperm, and therefore is only passed on in the female line. This, in fact, makes the above example of female lineages relevant for treating biological data. The other example is the Y-chromosome in humans and other mammals which is only carried by males (and determines the male gender) and therefore is only transmitted in the male lineage. The mathematically convenient solution, in contrast, is to

simply consider the smallest DNA segments, the single nucleotides. The biological problem here is that even though each nucleotide is derived from a unique parent, this usually cannot be identified from genetic data because in a given species, at most positions, most members share the same nucleotide. Nevertheless, for so-called SNPs, single nucleotide polymorphisms, consideration of single nucleotide positions can contain some useful population biological information. Even when we consider single nucleotide positions, however, we only get rid of the problem of recombination while absence of mutations, and of other processes of genetic rearrangement, then is still a hypothesis imposed. For simplicity of the model, we also consider the haploid case where each individual has only one set of genes. Thus, each such DNA segment in an individual is derived from one of the parents. (In the diploid case, each individual has two sets of genes. The genes corresponding to each other in those sets are selected from different parents, that is, one is taken from the mother and the complementary one from the father. This imposes additional restrictions when compared with the haploid case, but typically their effect is not so prominent.) For single nucleotides in the haploid case which we shall now consider the situation is formally the same as that one of the two parents of each individual is its mother. So, one might call that individual that gives the nucleotide in question the nucleotide parent for that particular position in the DNA. When in turn we consider only nucleotide parents, for a fixed position, then the situation is the same as before, with the only formal difference that two siblings can derive the nucleotide at that position from different parents. Therefore, the size of the population that has to be taken into account is  $2N$  in place of  $N$ .

In any case, for each such nucleotide position, we can perform the coalescent analysis and find the expected number of generations for having a single ancestor. Of course, the ancestors for the different positions will in general be different individuals. We can then also ask the following questions

1. What is the expected number of generations for finding a single ancestor for each position? That is, what is the expected maximal height of the coalescence trees for a given population?
2. In the corresponding ancestral population, how many individuals are ancestors for some position for the present sample? Those lucky ancestors then are ancestors of every individual in the present sample whereas the remaining members of the ancestral population then are not genetically represented at all in that present sample because for each position, there is only one ancestor by assumption.
3. Actually, the last issue is a little more subtle. In principle, an individual in the ancestral population can be a genealogical ancestor of a present one without being genetically represented in the latter. What are the chances for that? Here, one should essentially use some tree counting arguments. The pedigree graph contains many trees with a root in the ancestral population and leading down to the present population, and each

such root then represents a genealogical ancestor. Not all of these trees, however, arise from the coalescence processes just investigated.

So far, only some partial answers are known to these questions, see [17] for a brief discussion and references.



## Chapter 3

# Stochastic processes

### 3.1 Random variables

A general reference for this chapter is [15].

Let  $(\Omega, \Sigma, p)$  be a measure space, that is, a set  $\Omega$  with a probability measure  $p$  and a Sigma algebra  $\Sigma$  of measurable sets. A measurable function

$$X : \Omega \rightarrow \mathbb{R} \tag{3.1.1}$$

then is called a random variable. The possible values of  $X$  are called events. Thus, what is random here is not the function  $X$ , but rather its argument  $\omega \in \Omega$  that is considered to be drawn according to the measure  $p$ . The elements of  $\Omega$  may represent the possible outcomes of some experiment or observation. When tossing a coin, for example, there are two possible outcomes, heads  $H$  and tails  $T$ , and these then are the elements of  $\Omega$ . When the coin is tossed twice, the appropriate  $\Omega$  contains 4 elements,  $HH, HT, TH, TT$ . The random variable may be the number of heads; in the last example, it may take the values 0, 1, and 2. In this situation,  $X$  takes only discrete values, and whenever that is the case, we speak of a discrete random variable. When it takes its values in the non-negative integers  $\mathbb{N} = \{0, 1, 2, \dots\}$ , we have a counting variable. This will be the important case for us. For a discrete random variable, we have the (probability) mass function

$$p(x) := p(X = x) := p(\{\omega \in \Omega : X(\omega) = x\}), \tag{3.1.2}$$

with some abuse of notation (more precisely, we are using the symbol  $p$  both for the probability measure in  $\Omega$  and for the probability measure in  $\mathbb{R}$  induced by  $X$  – when there is a danger of confusion, we shall write  $p_X$  for the latter).

The distribution function of the random variable  $X$  is the function  $f(x) := p(X \leq x)$ .

When  $Y : \Omega \rightarrow \mathbb{R}$  is another random variable, we can consider the joint mass function

$$p(x, y) := p(X = x, Y = y) := p(\{\omega \in \Omega : X(\omega) = x, Y(\omega) = y\}) \tag{3.1.3}$$

and the conditional one,

$$p(y|x) := p(Y = y|X = x) = \frac{p(X = x, Y = y)}{p(X = x)} \quad (3.1.4)$$

whenever  $p(X = x) > 0$ , which satisfies

$$p(Y = y) = \sum_x p(Y = y|X = x)p(X = x) = \sum_x p(X = x, Y = y). \quad (3.1.5)$$

The random variables  $X$  and  $Y$  are called independent if the probabilities for all events  $x$  of  $X$  and  $y$  of  $Y$  satisfy

$$p(X = x, Y = y) = p(X = x)p(Y = y). \quad (3.1.6)$$

$X$  and  $Y$  are called identically distributed if for every  $z \in \mathbb{R}$

$$p(X = z) = p(Y = z). \quad (3.1.7)$$

The abbreviation “i.i.d” meaning independent and identically distributed is frequently used.

In the above example, we may take  $X$  as the number of heads. When the coin is fair, that is, when  $HH, HT, TH, TT$  each occur with probability  $1/4$ , we have  $p(0) = 1/4, p(1) = 1/2, p(2) = 1/4$ . More generally, when heads turn up with probability  $q$ , and if the results of the two tosses are independent of each other, we have  $p(0) = (1 - q)^2, p(1) = 2q(1 - q), p(2) = q^2$ . More generally, if we toss the coin  $n$  times, then  $\binom{n}{k}$  points in the corresponding  $\Omega$ , the set of the possible tossing sequences, yield  $k$  heads, and the mass function is

$$p(k) = \binom{n}{k} q^k (1 - q)^{n-k}. \quad (3.1.8)$$

This is the binomial distribution  $Bi(n, q)$ . When we let  $n \rightarrow \infty, q \rightarrow 0$  in such a manner that  $nq \rightarrow \lambda \neq 0$ , then, using the approximations  $\binom{n}{k} \approx \frac{n^k}{k!}$  and  $(1 - q)^{n-k} \approx (1 - q)^n = (1 - \frac{\lambda}{n})^n \approx e^{-\lambda}$ , we obtain the limit

$$\binom{n}{k} q^k (1 - q)^{n-k} \rightarrow \frac{\lambda^k}{k!} e^{-\lambda}. \quad (3.1.9)$$

This is the Poisson distribution  $Q(\lambda)$ .

The basic continuous distribution is the Gaussian distribution

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) \quad (3.1.10)$$

on  $\mathbb{R}$  with mean  $\mu$  and variance  $\sigma^2$ . Similarly, on  $\mathbb{R}^n$ , we have the multinomial Gaussian distribution

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right), \quad (3.1.11)$$

where  $\|\cdot\|$  denotes the Euclidean norm, and  $x, \mu \in \mathbb{R}^n$ .

**Definition 3.1.1.** The  $k$ -th moment of the discrete random variable  $X$  with mass function  $p$  is the expectation value of  $X^k$ ,

$$E(X^k) = \sum_x x^k p(x). \quad (3.1.12)$$

whenever that sum converges absolutely.

In fact, in order to make this consistent, one verifies more generally that for  $\phi : \mathbb{R} \rightarrow \mathbb{R}$

$$E(\phi(X)) = \sum_x \phi(x)p(x) \quad (3.1.13)$$

whenever the sum converges absolutely.

When the random variable  $X$  is not discrete, the above sums get replaced by an integral; for example

$$E(X^k) = \int_x x^k p(x) dx. \quad (3.1.14)$$

For the moment, however, we consider discrete random variables.

The first moment  $E(X)$  is of course the mean, average, or expectation value of  $X$ , and the second moment then yields its variance  $\text{var}(X) = E(X^2) - (E(X))^2 = E((X - E(X))^2)$ .

We can also consider conditional expectations and obtain the following lemma

**Lemma 3.1.1.** Let  $X$  and  $Y$  be discrete random variables on  $\Omega$ . The conditional expectation  $E(Y|X) =: \phi(X)$  satisfies

$$E(\phi(X)) = E(Y). \quad (3.1.15)$$

*Proof.* By (3.1.13) and (3.1.5)

$$\begin{aligned} E(\phi(X)) &= \sum_x \phi(x)p_X(x) = \sum_x \sum_y y p_{Y|X}(y|x)p_X(x) \\ &= \sum_{x,y} y p_Y(y) = E(Y). \end{aligned}$$

□

The random variables  $X$  and  $Y$  are called uncorrelated if

$$E(XY) = E(X)E(Y). \quad (3.1.16)$$

Independent random variables are uncorrelated, but not necessarily conversely. The next result is easy:

**Lemma 3.1.2.** For random variables  $X, Y$  and  $\alpha, \beta \in \mathbb{R}$ , we have

a)

$$E(\alpha X + \beta Y) = \alpha E(X) + \beta E(Y). \quad (3.1.17)$$

b)

$$\text{var}(\alpha X) = \alpha^2 \text{var}(X) \quad (3.1.18)$$

c) when  $X$  and  $Y$  are uncorrelated

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y). \quad (3.1.19)$$

□

a) simply expresses the linearity of the expectation value. Therefore, one may perform arbitrary *linear* operations with random variables without requiring that they be independent and have the corresponding operations on the expectation values.

We apply this lemma to the binomial distribution  $Bi(n, q)$  from (3.1.8), interpreted as  $n$  independent tossings of a coin. Since the value of the random variable  $H$ , the number of heads, in a single toss is 0 or 1, we have  $E(H^2) = E(H) = q$  and so the variance  $E(H^2) - (E(H))^2$  is  $q(1 - q)$ . Thus, the random variable for the event  $H$  has expectation value  $q$  and variance  $q(1 - q)$ . We therefore obtain from the lemma

$$E(Bi(n, q)) = nq, \quad \text{var}(Bi(n, q)) = nq(1 - q). \quad (3.1.20)$$

For the Poisson distribution  $Q(\lambda)$  (3.1.9) obtained above as a limit of binomial distributions, we obtain

$$E(Q(\lambda)) = \text{var}(Q(\lambda)) = \lambda. \quad (3.1.21)$$

There are various notions of convergence for a sequence  $(X_n)_{n \in \mathbb{N}}$  of random variables towards a random variable  $X$ .  $X_n$  converges to  $X$  almost surely if  $p(\omega \in \Omega : \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega)) = 1$ . It converges to  $X$  in  $k$ th mean ( $k \geq 1$ ) if all  $E(|X_n|^k) < \infty$  and  $\lim_{n \rightarrow \infty} E(|X_n - X|^k) = 0$ . The cases of most interest are  $k = 1$  and  $k = 2$ . For  $k = 2$ , one speaks of convergence in mean square.  $(X_n)$  converges to  $X$  in probability if for all  $\varepsilon > 0$ ,  $\lim_{n \rightarrow \infty} p(|X_n - X| > \varepsilon) = 0$  (as usual, the expression here is shorthand for  $p(\{\omega \in \Omega : |X_n(\omega) - X(\omega)| > \varepsilon\})$ ). It converges to  $X$  in distribution if  $\lim_{n \rightarrow \infty} p(X_n \leq x) = p(X \leq x)$  for all  $x \in \mathbb{R}$  for which the right hand side is continuous. Almost sure convergence and convergence in  $k$ th mean ( $k \geq 1$ ) each imply convergence in probability, and the latter in turn implies convergence in distribution, and convergence in  $k$ th mean implies convergence in  $l$ th mean for  $k > l \geq 1$ . There exist no other general implications between these notions of convergence.

We now state some fundamental convergence theorems, referring to [15] for proofs. The first is the **law of large numbers**:

**Theorem 3.1.1.** Let  $X_n, n \in \mathbb{N}$  be i.i.d random variables with  $E(|X_1|) < \infty$ . With  $\mu := E(X_1)$ , we then have

$$\frac{1}{n} \sum_{\nu=1}^n X_\nu \rightarrow \mu \quad \text{almost surely.} \quad (3.1.22)$$

If  $E(X_1^2) < \infty$ , the convergence takes also place in mean square.

The next is the **central limit theorem**:

**Theorem 3.1.2.** *Let  $(X_n)$  be a sequence of i.i.d. random variables with finite mean  $\mu$  and finite variance  $\sigma^2 \neq 0$ . Then, for  $S_n := \sum_{\nu=1}^n X_\nu$ , the distribution of*

$$\frac{S_n - n\mu}{\sqrt{n\sigma^2}} \text{ converges to the Gaussian distribution } \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right). \quad (3.1.23)$$

Let now

$$X : \Omega \rightarrow \mathbb{N} \quad (3.1.24)$$

be a discrete random variable that assumes only non-negative integer values.

**Definition 3.1.2.** The generating function of the random variable  $X$  is

$$G(s) := E(s^X) = \sum_{n=0}^{\infty} s^n p(n) \quad (3.1.25)$$

(defined for those values of  $s \in \mathbb{R}$  for which the sum converges).

Of course, the sequence  $p(n)$  can be recovered from the generating function by evaluating its  $k$ th derivative at 0:

$$G^{(k)}(0) = k! p(k) \quad (3.1.26)$$

As we shall now see, the derivatives of  $G$  at 1 also encode important properties of the sequence  $p(n)$ , namely its moments. Thus, when letting the argument  $s$  vary from 0 to 1, the generating function interpolates between the individual probabilities and collective properties of the distribution, the moments. Moreover, the generating function behaves in a very useful manner under composition of random processes and allows for a computation of moments of composed processes from the moments of the individual processes.

**Lemma 3.1.3.** *a) Let  $G$  be the generating function of  $X$ . Then*

$$E(X) = G'(1) \text{ and more generally } E(X(X-1)\cdots(X-k+1)) = G^{(k)}(1) \quad (3.1.27)$$

*whenever that  $k$ -th derivative of  $G$  at  $s = 1$  exists. (Thus, the moments of  $X$  can be computed recursively from the generating function.)*

*b) If  $X_1, \dots, X_N$  are independent random variables, their generating functions satisfy*

$$G_{\sum_{\nu=1}^N X_\nu}(s) = G_{X_1}(s) \cdots G_{X_N}(s). \quad (3.1.28)$$

*c) If  $X_1, X_2, \dots$  are independent and identically distributed random variables which then have the same generating function, denoted by  $G_X$ , and if  $N : \Omega \rightarrow \mathbb{N}$  is another random variable independent of the  $X_\nu$  with generating function  $G_N$ , then the random variable  $Y = X_1 + X_2 + \dots + X_N$  (that is, the number of random variables occurring in the sum is now a random variable itself) has the generating function*

$$G_Y(s) = G_N(G_X(s)). \quad (3.1.29)$$

*Proof.* a) is obvious, and b) follows from (3.1.16) applied to the independent random variables  $s^{X_\nu}$ . For c), from Lemma 3.1.1

$$\begin{aligned} G_Y(s) &= E(s^Y) = E(E(s^Y|N)) = \sum_n E(s^Y|N=n) p_N(n) \\ &= \sum_n E(s^{X_1+\dots+X_n}) p_N(n) \\ &= \sum_n E(s^{X_1}) \dots E(s^{X_n}) p_N(n) \text{ by b)} \\ &= \sum_n G_X(s)^n p_N(n) = G_N(G_X(s)). \end{aligned}$$

□

An alternative to the above polynomial generating function is the exponential one where we replace  $s$  in (3.1.25) by  $e^t$  to get

$$H(t) := E(e^{tX}) = \sum_{n=0}^{\infty} e^{nt} p(n). \quad (3.1.30)$$

Here, the moments of  $p$  can directly be computed from the derivatives of  $H(t)$  at  $t = 1$ . All the formal results that we demonstrate about  $G$  also hold for  $H$ . In fact, for most purposes, the generating function  $H$  is more convenient than  $G$ . In Theorem 3.4.1 below, however, we shall need a particular property of  $G$ , and this is the main reason why we are working here systematically with  $G$  in place of  $H$ . Also, it will be useful in our discussion of random graphs below.

A variant of  $H$  is the discrete Fourier transform, the so-called characteristic function

$$E(e^{itX}) = \sum_{n=0}^{\infty} e^{int} p(n) = \sum_n \sum_{\nu=0}^{\infty} \frac{(it)^\nu}{\nu!} n^\nu p(n), \quad (3.1.31)$$

also called the moment generating function, that similarly encodes the properties of the distribution  $p(n)$ .

## 3.2 Random processes

**Definition 3.2.1.** A random or stochastic process is a family  $X = (X_t)$  of random variables indexed by some set  $T \subset \mathbb{R}$ .

*Remark.* More generally, one can allow for the  $X_t$  to take values in some measurable space other than  $\mathbb{R}$ .

As  $X_t$  is a random variable, for each  $t$ , we get an induced probability distribution for the values of  $X$  by

$$p_t(S) := p(X_t \in S) := p(\{\omega \in \Omega : X_t(\omega) \in S\}) \text{ for a measurable } S \subset \mathbb{R}, \quad (3.2.32)$$

analogously to (3.1.2) (where we had only considered the case of discrete values). The random variables  $X_{t_1}$  and  $X_{t_2}$  (say  $t_1 < t_2$ ) are independent, see (3.1.6), if for  $S_1, S_2 \subset \mathbb{R}$

$$p(X_{t_1} \in S_1, X_{t_2} \in S_2) = p(X_{t_1} \in S_1) p(X_{t_2} \in S_2). \quad (3.2.33)$$

The process is called stationary if its finite dimensional distributions are time invariant, i.e.

$$p(X_{t_1+\tau} \in S_1, \dots, X_{t_n+\tau} \in S_n) = p(X_{t_1} \in S_1, \dots, X_{t_n} \in S_n) \quad (3.2.34)$$

for all  $t_1 < \dots < t_n$ ,  $S_1, \dots, S_n \subset \mathbb{R}$  and  $-\infty < \tau < \infty$ .

We shall now consider the case  $T = \mathbb{N}$ .

**Definition 3.2.2.** The random process  $X$  is called a Markov chain if

$$p(X_{n+1} = x | X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = p(X_{n+1} = x | X_n = x_n) \quad (3.2.35)$$

for all  $n \geq 1$ ,  $x, x_1, \dots, x_n \in \mathbb{R}$ .

**Definition 3.2.3.** The random process  $X$  is called a martingale if  $E(|X_n|) < \infty$  for all  $n$  and

$$E(X_{n+1} | X_1, X_2, \dots, X_n) = X_n. \quad (3.2.36)$$

We have the fundamental martingale convergence theorem

**Theorem 3.2.1.** A martingale  $X$  with  $E(X_n^2) < K < \infty$  for some  $K$  and all  $n$  converges to some random variable  $\Xi$  almost surely and in mean square.

### 3.3 Poisson processes and neural coding

In neurobiology, one studies spikes, that is, firings of neurons. Whereas there exist biophysical models for the generation of spikes on the basis of the electrochemical dynamics within neurons, see the Hodgkin-Huxley model described in 4.3.1 below, in more abstract models of information processing, spikes are conceptualized as discrete events occurring at time points, that is, events without temporal duration. This motivates

**Definition 3.3.1.** A stochastic process  $(N_t)_{t \in T}$  is called a *point process* when  $T$  is an interval in  $\mathbb{R}$  and the random variables  $N_t$  take only discrete values. It is called a *counting process* when  $N_0(\omega) = 0$  and  $N_t(\omega)$  counts the number of specified events in the interval  $[0, t]$ .

Of course, these concepts also apply to events other than spikes. Obviously, one can relax the normalizations, e.g., take some other  $t_0$  in place of 0, without much gain of generality.

Since the realization  $\omega \in \Omega$  will not play a significant role in the sequel, we shall omit it in our notation and henceforth write

$$N(t) \text{ in place of } N_t(\omega).$$

One may also consider the probability distribution of the time interval between events. The combination of these two points of view, that is, counting the number of events having occurred until time  $t$  vs. recording the temporal distance between subsequent events, will prove quite insightful. – If these time intervals are independently and identically distributed (a point to be returned to shortly), the process is called a *renewal process*. That is, the  $n$ th event occurs at time  $t_1 + t_2 + \dots + t_n$  where  $t_1, \dots, t_n$  are independent positive random variables that are identically distributed according to some probability density function  $p(t)$ . When this distribution does not depend on time, the process is called homogeneous, else inhomogeneous.

We return to the counting process.

**Definition 3.3.2.** A counting process  $N(t)$  is called locally continuous in probability if

$$\lim_{\epsilon \searrow 0} p(N(t + \epsilon) - N(t) \geq 1) = 0. \quad (3.3.37)$$

**Definition 3.3.3.** A counting process  $N(t)$  is said to have independent increments if the numbers of events in disjoint time intervals are independent.

**Definition 3.3.4.** A counting process  $N(t)$  is said to have stationary increments if the number of events in a time interval depends only on the length of that time interval.

**Definition 3.3.5.** A counting process  $N(t)$  that is locally continuous in probability and has independent and stationary increments is called a (homogeneous) Poisson process.

**Theorem 3.3.1.** For a Poisson process  $N(t)$

$$p(N(s + t) - N(s) = n) = e^{-rt} \frac{(rt)^n}{n!} \quad (3.3.38)$$

for  $n = 0, 1, \dots$ , for some constant  $r \geq 0$ .

Thus, the number of events produced by a Poisson process is distributed according to the Poisson distribution  $Q(rt)$ , see (3.1.9). The parameter  $rt$  of that distribution thus is proportional to the length  $t$  of the time interval, and the factor  $r$  is called the rate of the process.

*Proof.* We put

$$\rho(t) := p(N(t) - N(0) \geq 1) = p(N(s + t) - N(s) \geq 1) \quad (3.3.39)$$

where the last equality holds because  $N$  has stationary increments. We claim that

$$\rho(t) = 1 - e^{-rt} \quad (3.3.40)$$



for some constant  $r \geq 0$ . To see this, we start by observing that, by the assumption of independent increments, the probability  $\pi_0(t+s)$  that no event occurs between 0 and  $t+s$  equals the product of the probabilities that no event occurs in  $[0, t]$  and that no event occurs in  $[t, t+s]$ . By the assumption of stationary increments, the latter equals the probability that no event occurs in  $[0, s]$ , and hence

$$\pi_0(t+s) = \pi_0(t)\pi_0(s) \quad (3.3.41)$$

whence  $\pi_0(t) = \pi_0(0) \exp(-rt)$  for some constant  $r$ .  $r \geq 0$  since this probability is decreasing in  $t$ . Since  $\pi_0(0) = 1$  as  $N(0) = 0$ , we obtain

$$\pi_0(t) = e^{-rt}. \quad (3.3.42)$$

Since this is the probability that no event occurs until  $t$ , the probability  $\rho(t)$  that the first event occurs before  $t$  is  $1 - e^{-rt}$  which is (3.3.40).

Then

$$\rho(t) = rt + o(t) \text{ for } t \rightarrow 0. \quad (3.3.43)$$

Similarly, the probability  $\pi_1(t) := p(N(t) = 1)$  that precisely one event occurs between 0 and  $t+s$  satisfies

$$\pi_1(t+s) = \pi_0(t)\pi_1(s) + \pi_1(t)\pi_0(s) = e^{-rt}\pi_1(s) + e^{-rs}\pi_1(t). \quad (3.3.44)$$

Using (3.3.43), we obtain

$$\pi_1(t) = rte^{-rt}. \quad (3.3.45)$$

Iteratively, we obtain

$$\pi_n(t) := p(N(t) = n) = \frac{(rt)^n}{n!} e^{-rt}. \quad (3.3.46)$$

Recalling the assumption of stationary increments, this shows (3.3.38).  $\square$

As a consistency check, the reader should verify that a counting process  $N(t)$  given by (3.3.38) with independent increments conversely has stationary increments and is locally continuous in probability.

We also observe that

$$p(N(t) < \infty) = \sum_{n \geq 0} p(N(t) = n) = e^{-rt} \sum_{n \geq 0} \frac{(rt)^n}{n!} = 1, \quad (3.3.47)$$

and therefore, almost surely, only finitely many events take place in the finite interval  $[0, t]$ . Alternatively, this can also be directly be deduced from the assumptions. If there were infinitely many events in  $[0, t]$ , then there would also be infinitely many events in either  $[0, t/2]$  or  $[t/2, t]$ . By the assumption of stationary increments, then there would be infinitely many events in *both* these intervals. Repeating this argument, there would be infinitely many events in every subinterval of  $[0, t]$ . Consequently, every point would be an accumulation point of event, contradicting the assumption of local continuity in probability.

**Theorem 3.3.2.** *For a Poisson process with rate  $r$ , the time of the  $n$ th event is distributed according to*

$$p_n(t) = \frac{r(rt)^{n-1}}{(n-1)!} e^{-rt}. \quad (3.3.48)$$

*Proof.* We have

$$\frac{d}{dt} p(N(t) = n) = p_n(t) - p_{n+1}(t), \quad (3.3.49)$$

because the probability of having  $n$  events at time  $t$  goes up by the probability that the  $n$ th event occurs at time  $t$  and goes down with the probability that the  $(n+1)$ st event occurs then. Since  $p_0(t) = 0$  for  $t > 0$ , (3.3.48) follows iteratively from (3.3.46).  $\square$

The argument of the proof can also be reversed: The probability of observing the first event at  $t$  is given by the derivative of the probability for the first event occurring before  $t$ , that is, from (3.3.40), we obtain

$$p_1(t) = r e^{-rt}, \quad (3.3.50)$$

and as in the proof of Theorem 3.3.1, we can iterate that argument.

We also have the consistency relation

$$\sum_{\nu=0}^{n-1} \frac{(rt)^\nu}{\nu!} e^{-rt} = \int_t^\infty \frac{r(r\tau)^{n-1}}{(n-1)!} e^{-r\tau} d\tau \quad (3.3.51)$$

(which follows from repeated integration by parts). The lhs is the probability that at most  $n-1$  events have taken place up to time  $t$  (Theorem 3.3.1), and the rhs that the  $n$ th event occurs after that time (Theorem 3.3.2).

Assume that  $N(t) = n$ , that is, precisely  $n$  events occur in  $[0, t]$ . Since by the assumptions of stationary and independent increments, every  $n$ -tuple  $(t_1, \dots, t_n)$  of points in  $[0, t]$  has the same probability density of receiving those  $n$  events, that probability is given by

$$p(t_1, \dots, t_n) = \frac{1}{n!} r^n e^{-rt} \quad (3.3.52)$$

because this depends on the length  $t$  of the interval and yields (3.3.38) by integration w.r.t.  $t_1, \dots, t_n$  from 0 to  $t$ . The combinatorial factor  $\frac{1}{n!}$  arises, because the  $n$  events under consideration are indistinguishable. This also leads to the correct normalization, because when we integrate w.r.t.  $t_1, \dots, t_n$  and then sum over all  $n$ , we obtain 1, by the relation

$$\sum_n \frac{1}{n!} (rT)^n = e^{rT}. \quad (3.3.53)$$

**Theorem 3.3.3.** *Let  $N_1(t), N_2(t)$  be independent Poisson processes with rates  $r_1, r_2$ . Then the counting process  $N(t) = N_1(t) + N_2(t)$  is a Poisson process with rate  $r = r_1 + r_2$ .*

The *proof* is an obvious verification of the defining properties of a Poisson process.  $\square$

Likewise, we can multiply the rate of a Poisson process with a constant positive factor to obtain another Poisson process. In fact, when the factor is not constant, but is a function of time, we still obtain a process that is not significantly different as we shall now explore.

**Definition 3.3.6.** Let  $r(t)$  be a continuous positive function on  $\mathbb{R}_+$ , and put  $r_t := \int_0^t r(\tau)d\tau$ . A counting process  $N(t)$  with independent increments that satisfies

$$p(N(s+t) - N(s) = n) = e^{-r_t} \frac{(r_t)^n}{n!} \quad (3.3.54)$$

for  $n = 0, 1, \dots$  is called an (inhomogeneous) Poisson process with rate function  $r(t)$ .

Actually, this is not such a vast generalization (and we indicated that already by the typographical similarity between  $rt$  and  $r_t$ ). Since the parametrization of time is arbitrary, we have

**Theorem 3.3.4.** *Any inhomogeneous Poisson process can be transformed into a homogeneous one via a time reparametrization.*

*Proof.* Since  $r(t)$  is assumed positive,  $r_t$  is a strictly monotonically increasing, continuous function of  $t$  with  $r_0 = 0$  and  $\lim_{t \rightarrow \infty} r_t = \infty$ . It therefore has an inverse function  $\rho(t)$  with the same properties. We put

$$N'(t) := N(\rho(t)). \quad (3.3.55)$$

This counting process then is a homogeneous Poisson process with rate 1 since  $r_{\rho(t)} = t$ .  $\square$

The assumption that  $r(t)$  be strictly positive could be relaxed to nonpositivity. In that case, the time transformation would simply jump over those periods where  $r(t)$  vanishes.

An equivalent characterisation of an inhomogeneous Poisson process with rate function  $r(t)$  is that it be a counting process with independent increments, satisfying

$$p(N(t+\epsilon) - N(t) = 1) = r(\epsilon) + o(\epsilon) \text{ and } p(N(t+\epsilon) - N(t) \geq 2) = o(\epsilon) \text{ for } \epsilon \rightarrow 0. \quad (3.3.56)$$

The first relation is a quantitative version of the local continuity in probability. The second relation automatically holds in the homogeneous case, but one needs to assume this in the inhomogeneous case.

For an inhomogeneous Poisson process, the probability density for observing events at the times  $t_1, \dots, t_n$  is

$$p(t_1, \dots, t_n) = \frac{1}{n!} \exp\left(-\int_0^t r(\tau)d\tau\right) \prod_{i=1}^n r(t_i). \quad (3.3.57)$$

Using the time reparametrization of Theorem 3.3.4, this is deduced from the formula (3.3.52) for the homogeneous case. Upon integration, it yields (3.3.54), analogously to the homogeneous case.

From (3.3.54), the expected number of events in the interval  $[0, t]$  is

$$\sum_n np_t(n) = \frac{(r_t)^n}{(n-1)!} e^{-r_t} = r_t = \int_0^t r(\tau) d\tau \quad (3.3.58)$$

(see (3.3.53)).

Inhomogeneous Poisson processes are important models in theoretical neurobiology (see e.g. [8]) because the rate  $r(t)$  that represents the spiking or firing rate of a neuron can then depend on the stimulus  $S$ , that is, on the input received by the neuron. Thus,

$$r(t) = r(t; S), \quad (3.3.59)$$

and this represents the coding scheme of the neuron under consideration. In other words, when receiving the stimulus  $S$ , the neuronal firing follows a Poisson process with rate  $r(t; S)$ . For simplicity, we may assume that the stimulus is received at time  $t = 0$ . A basic model assumes that the neuron has a preferred or optimal stimulus  $S_0$ , and that the stimulus is translated via a Gaussian tuning function into the firing rate

$$r(t; S) = c \exp\left(-\frac{d^2(S, S_0)}{2\sigma^2}\right). \quad (3.3.60)$$

Here,  $d(., .)$  is some metric in the input space, for example a Euclidean one, that is,  $d(S, S_0) = \|S - S_0\|$ ;  $c$  and  $\sigma^2$  (variance) are parameters. Obviously, other coding schemes are possible.

From (3.3.57), we see that the probability distribution for observing spikes precisely at the times  $t_1, \dots, t_n$  in the interval  $[0, t]$ , given the input  $S$  is

$$p(t_1, \dots, t_n | S) = \frac{1}{n!} \exp\left(-\int_0^t r(\tau; S) d\tau\right) \prod_{i=1}^n r(t_i; S). \quad (3.3.61)$$

Bayes' formula then yields the fundamental relationship for decoding the spike train produced by the neuron, that is, an estimate for the distribution of signals contingent upon the recorded spike train  $t_1, \dots, t_n$ ,

$$p(S | t_1, \dots, t_n) = p(t_1, \dots, t_n | S) \frac{p(S)}{p(t_1, \dots, t_n)}. \quad (3.3.62)$$

Here,  $p(S)$  is a prior estimate for the distribution of stimuli (that may have been obtained as the result of some learning process).  $p(t_1, \dots, t_n)$  simply represents some normalization factor.

Of course, (3.3.61) can also be used for other estimation schemes. For example, maximum likelihood selects a stimulus  $\bar{S}$  that has caused an observed spike sequence with the highest probability, that is,  $\bar{S} = \operatorname{argmax} p(t_1, \dots, t_n | S)$  in (3.3.61).

Clearly, a biological neuron does not operate according to a Poisson process. A spike is not an instantaneous event, but the generation of an action potential has a positive, although rather short, duration. Also, two spikes cannot be fired in too rapid succession because after a spike is fired, a neuron goes through a refractory period before it can generate the next spike. Thus, if one wants to include these aspects, one needs a biophysical in place of a purely phenomenological model. Such models exist, and we shall introduce and investigate below the basic one, the Hodgkin-Huxley model. Notwithstanding its lack of biophysical realism, however, Poisson type models are very important in the neurosciences because, on one hand, they relate well to the experimental practice of recording spikes, and on the other hand, they can be the basis for models of information transmission in neural systems.

### 3.4 Branching processes

References for this section are [23, 16].

We start with the simplest branching process, the Galton-Watson process. Here, each individual lives in a fixed generation  $n$  and independently of all other individuals produces a random number of offspring that become members of generation  $n + 1$ . This random variable, the number of offspring, is the same for all individuals in all generations. Thus, the numbers of offspring for the individuals are independent and identically distributed random variables. We denote their common generating function by  $G(s)$ . We also assume that there is a positive probability for having more than one offspring. If the probability of having  $m$  offspring is  $p(m)$ , this means that  $p(0) + p(1) < 1$ .

Let the random variable  $Z_n$  denote the size of generation  $n$ . One usually assumes that the process starts with a single individual in generation 0, that is,  $Z_0 = 1$ . Let  $G_n(s) = E(s^{Z_n})$  be the generating function of  $Z_n$ .

**Lemma 3.4.1.**  $G_n$  is the  $n$ -th fold iterate of  $G$ ,

$$G_n(s) = G \circ \dots \circ G(s), \quad (3.4.63)$$

and thus also for  $m, n \in \mathbb{N}$

$$G_{m+n}(s) = G_m(G_n(s)). \quad (3.4.64)$$

*Proof.* We shall show (3.4.64) which easily implies (3.4.63) by iteration. Let the random variable  $Y_i$  denote the number of members of the  $(m + n)$ th generation that derive from member  $i$  of the  $m$ th one. We then have

$$Z_{m+n} = Y_1 + \dots + Y_{Z_m}. \quad (3.4.65)$$

By our assumptions, the  $Y_i$  are independent and identically distributed, in fact identical to  $Z_n$ , the number of offspring deriving from an individual  $n$  generations ago. Lemma 3.1.3 c) then yields the claim.  $\square$

**Corollary 3.4.1.** *Let  $\mu := E(Z_1)$  and  $\sigma^2 := \text{var}(Z_1)$ . Then*

$$E(Z_n) = \mu^n \text{ and } \text{var}(Z_n) = \begin{cases} n\sigma^2 & \text{if } \mu = 1 \\ \frac{\sigma^2(\mu^n - 1)\mu^{n-1}}{\mu - 1} & \text{if } \mu \neq 1. \end{cases} \quad (3.4.66)$$

*Proof.* Differentiating  $G_n(s) = G(G_{n-1}(s))$  at  $s = 1$  yields

$$E(Z_n) = \mu E(Z_{n-1}) \quad (3.4.67)$$

from which the first equation follows by iteration. Differentiating twice gives  $G_n''(1) = G''(1)G_{n-1}'(1)^2 + G'(1)G_{n-1}''(1)$  which yields the second equation.  $\square$

In view of this result, we call the process subcritical, critical, supercritical when  $\mu < 1, = 1, > 1$ , resp. In the sub- (super-)critical, we thus expect the population to shrink (grow) while in the critical it is expected to stay the same. This might lead one to expect that the population will continue forever, but that is not true as we shall now find out from asking whether the population will eventually become extinct, that is

$$Z_n = 0 \text{ for some } n \in \mathbb{N}, \text{ and then of course also for all } m \geq n. \quad (3.4.68)$$

Some observations are obvious:

- If  $p(0) = 0$ , that is, if every individual always has at least one offspring, then the population cannot become extinct. Therefore, we shall assume now

$$p(0) > 0. \quad (3.4.69)$$

- When  $p(\nu) = 0$  for  $\nu \geq 2$ , then  $p(0) + p(1) = 1$ . We have excluded  $p(0) = 0$  and therefore cannot have the trivial case  $p(1) = 1$ , that is, that every individual has precisely one offspring so that the population size will always remain constant. Consequently, the population should also become extinct because every individual then has either no or one offspring, and the former with a positive probability. Therefore, the population will certainly decrease. Therefore, we shall assume now

$$p(0) + p(1) < 1. \quad (3.4.70)$$

**Theorem 3.4.1.** *The extinction probability  $q_{\text{ext}}$  of the process, that is, the probability that for some  $n \in \mathbb{N}$  we have  $Z_n = 0$ , equals the smallest root of the equation  $G(s) = s$ . For  $\mu \leq 1$ , we have  $q_{\text{ext}} = 1$ , that is, the population becomes extinct almost surely, while for  $\mu > 1$ , we have  $q_{\text{ext}} < 1$ , that is, the population has a positive probability of surviving forever.*

*Proof.* We observe that  $G(s)$  as a power series with non-negative coefficients  $p(\nu)$  and  $p(0) + p(1) < 1$  by our initial assumptions is increasing and strictly convex for  $s \in [0, 1]$  and satisfies  $G(0) = p(0), G(1) = 1$ . When  $\mu = G'(1) \leq 1$ , then  $G(s) > s$  for  $s \in [0, 1)$ , while for  $\mu > 1$ ,  $G(s) = s$  has a unique root in  $[0, 1)$ . These properties implies that  $G(s) = s$  has a smallest root which we denote by  $q$ , and  $q = 1$  for  $\mu \leq 1$ , but  $q < 1$  for  $\mu > 1$ . Moreover, for  $s \in [0, q)$ , in particular for  $s = 0$ , the iterates  $G_n(s) = G \circ \dots \circ G(s)$  increase monotonically towards  $q$  while for  $s \in (q, 1)$  it decreases monotonically towards  $q$  for  $n \rightarrow \infty$ . (In terms of dynamical systems, this simply expresses the stability of the fixed point of  $G$  at  $q$  under dynamical iteration of  $G$ , which always holds when the graph of a function  $G$  intersects the diagonal from above at a fixed point.) We recall from Lemma 3.4.1 that  $G_n$  is the generating function for  $Z_n$ . Thus

$$\begin{aligned} q &= \lim_{n \rightarrow \infty} G_n(0) = \lim p(Z_n = 0) = \lim p(Z_\nu = 0 \text{ for some } \nu \leq n) \\ &= p(Z_\nu = 0 \text{ for some } \nu \in \mathbb{N}) = p(\lim Z_n = 0) \end{aligned}$$

is the extinction probability.  $\square$

Thus, we see that due to the fluctuations in the number of offspring, a finite population may become extinct in finite time. It will do so almost surely when the expected number of offspring is at most 1 – even when it is 1 –, and it will also go extinct with a positive probability when that number is larger than 1. One may consider this as a finite size effect, in the sense that when we go to the limit of large populations (under appropriate technical conditions), the random fluctuations will average out and the population will expand or shrink deterministically at the rate  $\mu$ .

A Galton-Watson branching process  $(Z_n)$  is a Markov process by (3.4.67). The normalized process  $W_n := \frac{Z_n}{E(Z_n)}$  then is a martingale by (3.4.66),

$$E(W_{n+1} | W_1, \dots, W_n) = W_n. \quad (3.4.71)$$

We now look at the situation where the expectation values  $\mu(n)$  for the number of offspring of an individual in generation  $n$  vary, i.e.,  $\mu$  is a random variable itself (defined on  $\mathbb{N}$ ). This is called a branching process in a random environment. The population then grows from time 0 to time  $n$  by the sequence  $\mu(0), \mu(1), \dots, \mu(n-1)$  which is equivalent to growth by the geometric mean  $(\mu(0) \cdots \mu(n-1))^{1/n}$  in one step. In order to apply the law of large numbers, we need to convert this product into a sum,

$$(\mu(0) \cdots \mu(n-1))^{1/n} = \exp\left(\frac{1}{n}(\log \mu(0) + \cdots + \log \mu(n-1))\right) \quad (3.4.72)$$

and conclude by the law of large numbers

$$\lim_{n \rightarrow \infty} \frac{1}{n}(\log \mu(0) + \cdots + \log \mu(n-1)) = E(\log \mu) \quad (3.4.73)$$

with probability 1. Thus, the asymptotic population growth rate is

$$\lim_{n \rightarrow \infty} (\mu(0) \cdots \mu(n-1))^{1/n} = \exp(E(\log \mu)) \quad (3.4.74)$$

which may be smaller than  $E(\mu)$ . In particular, even when the latter may be  $> 1$ , the process may still be subcritical, that is become extinct with positive probability, because of the fluctuations in the environment. So, once more we see that when a finite population is subjected to random effects its extinction probability may increase even when the expected growth rate stays the same.

The Galton-Watson process is the simplest branching process, and many generalizations are possible. One of them is to allow for individuals of different types  $j = 1, \dots, m$ . For each type, the distribution of the types of its offspring may be different. We then consider the matrix  $M = (m_{ij})$  where  $m_{ij}$  is the expected number of offspring of type  $j$  of an individual of type  $i$ .<sup>1</sup> All entries of  $M$  are non-negative. The expectation value of the number  $Z_{j,n}$  of individuals of type  $j$  in generation  $n$  is then

$$E(Z_{j,n}) = \sum_{i=1}^m E(Z_{i,n-1})m_{ij} \quad (3.4.75)$$

by linearity. In vector notation, with  $E(Z_n) = (E(Z_{1,n}), \dots, E(Z_{m,n}))^T$  ( $T$  denoting transpose),

$$E(Z_n)^T = E(Z_{n-1})^T M = E(Z_0)^T M^n. \quad (3.4.76)$$

Of course, when as before, we specify the initial population  $Z_0$ , we can drop the last expectation  $E$  to get  $E(Z_n)^T = Z_0^T M^n$ .

We shall now apply the theory of Perron-Frobenius to the non-negative matrix  $M$ . That theory is summarized in

**Lemma 3.4.2.** *Let  $M$  be an  $m \times m$  matrix with non-negative entries which is irreducible in the sense that all the entries of  $M^\nu$  are even positive for some  $\nu \in \mathbb{N}$ . Then  $M$  has a simple eigenvalue  $\rho$  that is real and positive and larger than the absolute value of any other eigenvalue. It possesses a left eigenvector  $u = (u^1, \dots, u^m)^T$  (that is,  $u^T M = \rho u^T$ , or in components,  $\sum_j u^j m_{ji} = \rho u^i$ ) and a right eigenvector  $v = (v^1, \dots, v^m)^t$  (i.e.,  $Mv = \rho v$ ) that both have positive entries. We can normalize them by*

$$\sum_{j=1}^m u^j = 1, \quad \sum_{j=1}^m v^j = 1. \quad (3.4.77)$$

With  $M_0 := (v^i u^j)_{i,j=1,\dots,m}$ , we then have

$$M^n = \rho^n M_0 + \tilde{M}^n, \quad (3.4.78)$$

---

<sup>1</sup>Later on, we shall consider the probability  $p^i(n_1, \dots, n_m)$  that an individual of type  $i$  produces  $n_k$  offspring of type  $k$ . Then  $m_{ij} = \sum_{n_1, \dots, n_m=1}^{\infty} n_j p^i(n_1, \dots, n_m)$ .



with  $M_0\tilde{M} = \tilde{M}M_0 = 0$  and  $|\tilde{M}^n| \leq \text{const } \tilde{\rho}^n$  for some  $\tilde{\rho} < \rho$ . In particular, for any non-trivial  $w$  with non-negative entries, the iterates  $w^T M^n$  grow like  $\rho^n$  in norm.

Returning to (3.4.76), we see that when  $E(Z_0) = cu$ , i.e., is a multiple of the left eigenvector  $u$  for the maximal eigenvalue  $\rho$ , we obtain

$$E(Z_n)^T = cu^T M^n = cu^T \rho^n. \quad (3.4.79)$$

We conclude that the process is subcritical, critical or supercritical depending on whether  $\rho < 1, = 1, > 1$ . In fact, there seems to be a small caveat here, namely that for (3.4.79), we had assumed that the initial composition  $Z_0$  of the population is a multiple of the positive eigenvector  $u$ . In order to be able to drop that assumption, we now assume that the matrix  $M$  is irreducible as in the Perron-Frobenius theorem. Such a process is called indecomposable. Then the value of  $\rho$  determines the asymptotic behavior of the population size for any initial configuration of the population.

We can also set up the generating function formalism as before, with the sole difference that all expressions now become vectors in place of scalars. The generating vector is  $G = (G^1, \dots, G^m)$  with

$$G^j(s_1, \dots, s_m) = \sum_{n_1, \dots, n_m} s_1^{n_1} \cdots s_m^{n_m} p^j(n_1, \dots, n_m) \quad (3.4.80)$$

where  $p^j(n_1, \dots, n_m)$  is the probability that an individual of type  $j$  produces  $n_i$  offspring of type  $i$ , for  $i = 1, \dots, m$ . We also have a vector  $q = (q_1, \dots, q_m)$  where  $q_j$  is the extinction probability for a population starting with a single individual of type  $j$ . Then, as in Theorem 3.4.1, the vector  $q$  is determined as the componentwise smallest root of the vector fixed point equation

$$q = G(q). \quad (3.4.81)$$

With some simple tricks, many different processes can be captured by multi-type Galton-Watson processes:

1. Given a single-type Galton-Watson process, we want to know the total number of individuals up to time  $n$ . We then simply define a second type of individual in the original process, the dead type. Type 1 corresponds to the original one, and it produces offspring of that type 1 according to the original rule, plus 1 individual of type 2, that is, it dies, as already assumed in the original process. An individual of type 2 produces one offspring of type 2, that is, it stays dead. By this token, the individuals from previous generations remain in all future generations, simply as corpses, i.e. as type 2 individuals. The transition matrix then is

$$\begin{pmatrix} \mu & 1 \\ 0 & 1 \end{pmatrix},$$

with  $\mu = E(Z_1)$  as above.

2. We consider a population consisting of two sexes. Females produce offspring of either type, with equal probabilities, while males do not reproduce. Our transition matrix then is of the form

$$\begin{pmatrix} \mu/2 & \mu/2 \\ 0 & 0 \end{pmatrix}.$$

While this process again is decomposable, here the two types still grow or shrink at the same rate. In fact, however, this is only an incomplete model of populations with sexual reproduction because they can become extinct not only for the reason that the population size goes to 0, but also when one of the two sexes disappears. That aspect needs to be modeled separately through the choice of a mating function, that is, by a rule how the number of offspring depends on the numbers of the two sexes in the population.

3. One can also include dependencies between siblings. That means that the expected numbers and types of offspring of an individual depend not only on its own type, but also on that of its siblings. While this violates the independence hypothesis in Cor.3.4.1, it turns out that dependencies in the same generation do not affect the expected population size. Again, this can be seen through a simple trick, namely by formally considering the sibship (brood, litter) as the individuals in the process. Different such sibships then produce different sibships in the next generation.

We consider the case of altruistic siblings, for example where the eldest one may forego its own offspring for helping his younger siblings to raise additional offspring. Of course, the Galton-Watson assumptions that reproduction takes place at discrete time steps, and that each individual can reproduce only at age 1 make the distinction between older and younger siblings impossible if taken literally, but for the sake of the argument we assume that some litters contain an altruistic member – which we then simply label as the “eldest” – whereas others don’t. In other words, we have two types of litters, one with an altruistic member, and the other one without. As before,  $m_{ij}$  is the expected number of progeny of type  $j$  produced by a litter of type  $i$ . By the Perron-Frobenius theorem, if  $M$  is irreducible, it has an eigenvalue  $\rho$  of largest absolute value that is positive and simple.  $\rho$  describes the growth rate, the corresponding normalized left eigenvector  $u$  yields the relative asymptotic contributions of the types to future generations while the right eigenvector  $v$  describes the asymptotic distribution of the process in case  $\rho > 1$  in the following sense: We assume that the process is regular, i.e. that the probability for a litter producing more than one progeny is positive, and that the second moments of progeny distribution are finite. Then by the martingale convergence theorem 3.2.1, see (3.4.71),  $Z_n/\rho^n$  converges to a vector  $w$  which (with probability 1) is a positive (except in obvious trivial cases) multiple of  $v$ . In that case, the process will not become extinct with positive probability while in case  $\rho \leq 1$ , it goes extinct with probability 1 as is typical for

branching processes. In fact, in our simple situation with two types, that eigenvalue is given by

$$\rho = \frac{1}{2}(m_{11} + m_{22}) + \sqrt{\frac{1}{4}(m_{11} - m_{22})^2 + m_{12}m_{21}}. \quad (3.4.82)$$

Let us assume that type 1 is the altruistic, and type 2 the non-altruistic one, and that  $m_{21} = 0$ , i.e., that a non-altruistic cannot produce an altruistic one ( $M$  then is no longer irreducible but the needed results from the Perron-Frobenius theorem still hold here). Then, if  $m_{11} > m_{22}$ , i.e. the altruistic litters reproduce more successfully than the other ones,

$$\rho = m_{11}, \quad (3.4.83)$$

and the corresponding eigenvector is proportional to

$$\begin{pmatrix} m_{11} - m_{22} \\ m_{12} \end{pmatrix}. \quad (3.4.84)$$

Thus, the non-altruists only survive as a nontrivial fraction of the total population<sup>2</sup> if  $m_{12} > 0$ , that is if they are also produced by the altruists. Of course, this is rather obvious.

If altruism is caused by a single gene, then one can use the methods of mathematical population genetics to compute the transmission probability of the responsible allele in a sexually reproducing population. The point of our example is that the altruistic allele has to be present in the sibling labeled “eldest” for being effective, but that it can only be transmitted to the next generation when also present in his siblings for whose own behavior the allele is irrelevant. In this manner, the coefficient  $m_{12}$  can be determined.

### 3.5 Random graphs

Equipped with tools from stochastic analysis, we now return to graph theory and discuss stochastic constructions of graphs. A good reference that we shall partly follow is [29].

The idea of Erdős and Rényi[10] that started the whole field was to not specify a graph explicitly, but rather only its generic type by selecting edges between nodes randomly, depending on a single parameter, the edge probability  $p$ . In a random graph, for any pair of nodes, there is thus an edge between them with probability  $p$ . If the network has  $N$  nodes, then each node has  $N - 1$  possible recipients for an edge. Thus, the average degree of a node is

$$z := (N - 1)p. \quad (3.5.85)$$

---

<sup>2</sup>Here, we are considering the population of litters, and not of individuals. For the latter, one would need to multiply these coefficients by the litter sizes.

Moreover, the probability that a given node has degree  $k$  in an Erdős-Rényi graph is

$$p_k = \binom{N-1}{k} p^k (1-p)^{N-1-k} \quad (3.5.86)$$

because the degree happens to be  $k$  when precisely  $p$  out of the  $N-1$  possible edges from the given node are chosen, and each of them is chosen with probability  $p$  and not chosen with probability  $1-p$ . Thus, the degree distribution is binomial, and for  $N \gg kz$ , this is approximated by the Poisson distribution

$$p_k = \frac{z^k e^{-z}}{k!}. \quad (3.5.87)$$

(and so  $z = \langle k \rangle = \sum_k k p_k$  (cf. (3.1.9), (3.1.21) on the Poisson distribution and (3.1.27) on the generating function).)

For an Erdős-Rényi graph, one can also compute the distribution of the number of second neighbors of a given node, that is, the number of neighbors of its neighbors, discarding of course the original node itself as well as all its direct neighbors that also happen to be connected with another neighbor. However, since there is no tendency to clustering in the construction, the probability that a second neighbor is also a first neighbor behaves like  $1/N$  and so becomes negligible for large  $N$ . Now, however, the degree distribution of first order neighbors of some node is different from the degree distribution of all the nodes in the random graph, because the probability that an edge leads to a particular node is proportional to that node's degree so that a node of degree  $k$  has a  $k$ -fold increased chance of receiving an edge. Therefore, the probability distribution of our first neighbors is proportional to  $k p_k$ , that is, given by  $\frac{k p_k}{\sum_l l p_l}$ , instead of  $p_k$ , the one for all the nodes in the graph. Since such a first neighbor of degree  $k$  has  $k-1$  edges leading away from the original node, the distribution for having  $k$  second neighbors via one particular one of its neighbors is then given, after shifting the index by 1, by

$$q_k = \frac{(k+1)p_{k+1}}{\sum_l l p_l}. \quad (3.5.88)$$

Thus, to obtain the number of second neighbors, we need to sum over the first neighbors, since, as argued, we can neglect clustering in this model. Thus, the mean number of second neighbors is obtained by multiplying the expected number of second neighbors via a particular first neighbor, that is,  $\sum k q_k$ , by the expected number of first neighbors,  $z = \sum k p_k$ . So, we obtain for that number

$$\sum_l l p_l \sum_k k q_k = \sum_{k=0}^{\infty} k(k+1)p_{k+1} = \sum_{k=0}^{\infty} (k-1)k p_k = \langle k^2 \rangle - \langle k \rangle. \quad (3.5.89)$$

We recall from 3.1 that such probability distributions can be encoded in probability generating functions (see (3.1.25)). If we have a probability distribution

$p_k$  as above on the non-negative integers, we have the generating function (cf. (3.1.25))

$$G_p(x) := \sum_{k=0}^{\infty} p_k x^k. \quad (3.5.90)$$

Likewise, the above distribution for the number of second neighbors then is encoded by

$$G_q(x) = \sum_{k=0}^{\infty} q_k x^k = \frac{\sum_k (k+1) p_{k+1} x^k}{\sum_l l p_l} = \frac{G'_p(x)}{z}. \quad (3.5.91)$$

When we insert the Poisson distribution (3.5.87), we obtain

$$G_p(x) = e^{-z} \sum_{k=0}^{\infty} \frac{z^k}{k!} x^k = e^{z(x-1)} \quad (3.5.92)$$

and from (3.5.91) then also

$$G_q(x) = e^{z(x-1)} \quad (3.5.93)$$

Thus, for an Erdős-Rényi graph, the two generating functions agree. This is quite useful for deriving analytical results.

When we construct a graph by a stochastic process like the one of Erdős-Rényi, the resulting structure need not be connected, but may have several components. In Chapter 2, it was part of the definition of a graph to be connected, but we drop that now because it will complicate our discussion of random graphs. Thus, an Erdős-Rényi graph may have several connected components. In order to understand this better, one lets  $N$  tend to  $\infty$  while keeping  $z = (N-1)p \sim Np$  from (3.5.85) fixed. It will then depend on the value of that parameter  $z$  whether the graph can be expected to contain a giant component or not. Here, a giant component is one that contains a positive fraction of the number of all vertices in the graph. More precisely, when  $z$  is above a critical threshold, we expect that our graph contains a component with at least  $\delta N$  vertices, for some  $\delta > 0$  that does not depend on  $N$ . Below that critical value, all components should have an average size that stays bounded as  $N \rightarrow \infty$ . This fact, and the computation of the value of  $z$  where that phase transition occurs, are the basic results of the theory of random graphs. Following [29], we now present a self-consistency argument to derive those results. (For a mathematically more rigorous treatment, we refer to [3].) We already noted that the clustering coefficients will tend to 0 for  $N \rightarrow \infty$ . Therefore, we expect components of bounded size not to contain any triangles. With the same kind of heuristic reasonings, we even assume that all our finite components do not contain cycles, that is, are trees. – Suppose we then randomly choose an edge in our graph, take one of its ends  $i_0$  and look at all the nodes that can be reached via other edges from  $i_0$ . Let the number of those nodes be  $x$ . We can then consider the generating function  $H_1(x)$  for the distribution  $p(x)$ . When we go from  $i_0$  to any of its new neighbors, that is,

other neighbors than the one from our edge that we started with, we are in the same situation as before, and for any of those neighbors we can again look at the number of nodes that can be reached from it via edges other than the one we arrived at it from  $i_0$ . That number  $x$  again is distributed according to the generating function  $H_1(x)$ . This then directly leads to a self-consistency equation. Namely, the above number  $k$  of new neighbors of  $i_0$  is distributed according to  $q_k$  from (3.5.88). Also, we recall from (3.1.28) that the generating function for a sum of independent processes is the product of the individual generating functions. Thus, we need to take the product of  $k$  factors  $H_1$ , one for each new neighbor of  $i_0$  weighted with the probabilities  $q_k$  and one additional factor  $x$  to account for  $i_0$  itself<sup>3</sup>,

$$H_1(x) = x \sum_{k=0}^{\infty} q_k (H_1(x))^k = x G_q(H_1(x)) \quad (3.5.94)$$

by (3.5.91). From this fixed point equation, we can compute  $H_1(x)$ . We can then easily determine the distribution  $H_0(x)$  for the total size of a finite component. Namely, take any vertex  $i_1$ ; it has  $k$  neighbors with probability  $p_k$ , and from each of them, we expect to reach a number  $x$  of further vertices distributed according to  $H_1(x)$ . Therefore, by the same reasoning as before,

$$H_0(x) = x \sum_{k=0}^{\infty} p_k (H_1(x))^k = x G_p(H_1(x)). \quad (3.5.95)$$

All generating functions  $G(x) = \sum_k x^k p(k)$  satisfy  $G(1) = 1$ , see (3.1.25); in particular, this holds for  $G_p$  and  $H_1$ . Likewise, the expectation value  $\langle k \rangle$  of  $k$  is given by  $G'(1)$ , see (3.1.27). Therefore, when all components are finite, the mean component size is

$$\langle x \rangle = H_0'(1) = 1 + G_p'(1) H_1'(1) = 1 + \frac{G_p'(1)}{1 - G_q'(1)} \quad (3.5.96)$$

with the help of (3.5.94). This becomes infinite when  $G_q'(1)$  approaches 1, and this then is the phase transition where a giant component appears. From (3.5.91), one can then compute the phase transition value. Above that transition value, the former analysis is no longer valid, but it can nevertheless be used to derive useful results. The point is that while for a giant component we can no longer neglect clustering effects it still applies to the finite components. When  $\sigma$  is the fraction of nodes in the giant component, we then have, using (3.5.95)

$$\sigma = 1 - H_0(1) = 1 - G_p(s), \quad (3.5.97)$$

with  $s = H_1(1)$  solving, by (3.5.94),

$$s = G_q(s) \quad (3.5.98)$$

---

<sup>3</sup>We have the single node  $i_0$  with probability 1, and so, the generating function is simply  $1 \cdot x$ .

(as in Theorem 3.4.1).

There exist other methods to reduce this percolation analysis to phase transition models in statistical mechanics, and we now discuss one such approach. A random graph  $\Gamma$  is a member of an ensemble defined by the parameters  $N$  and  $p$ . Its probability in this ensemble is given by

$$P_\Gamma = p^{l(\Gamma)}(1-p)^{\binom{N}{2}-l(\Gamma)} = \exp\left(-\frac{pN^2}{2} + pN\left(\frac{1}{2} - \frac{pN}{4} + \frac{l(\Gamma)}{N} + o(1)\right)\right)p^{l(\Gamma)}, \quad (3.5.99)$$

$l(\Gamma)$  denoting the number of edges. The probability of a random graph to have  $m$  components then is written as

$$P_m = \sum_{\Gamma} P_\Gamma \delta(m, m(\Gamma)), \quad (3.5.100)$$

$\delta(m, m(\Gamma))$  being the Kronecker delta. In order to study the distribution of components, one considers

$$P_\Gamma(q) := \frac{1}{z(q)} P_\Gamma q^{m(\Gamma)} \quad (3.5.101)$$

with the normalizing factor

$$z(q) := \sum_{\Gamma} P_\Gamma q^{m(\Gamma)} = \sum_m P_m q^m. \quad (3.5.102)$$

This quantity is related to the properties of a model from statistical mechanics, the Ising model with  $q$  states, also called the Potts model. In this model, one has  $N$  spin variables  $\sigma_i$  that can take one of  $q$  distinct values  $\sigma = 0, 1, \dots, q-1$ . The energy function of this model (in the so-called mean field variant) is

$$E(\{\sigma_i\}) := -\frac{1}{N} \sum_{i<j} \delta(\sigma_i, \sigma_j) - h \sum_{\sigma=0}^{q-1} f_\sigma \sum_i \delta(\sigma_i, \sigma) \quad (3.5.103)$$

where  $f_\sigma$  is an external field in the spin direction  $\sigma$ , multiplied with the strength  $h$ . One introduces a so-called inverse temperature  $\beta$  and encodes the thermodynamic properties of the model at  $\beta$  in the partition function

$$Z_\beta(q) := \sum_{\{\sigma_i\}} \exp(-\beta E(\{\sigma_i\})); \quad (3.5.104)$$

the sum is taken over all possible spin configurations  $\{\sigma_i\}$ . The partition function can be rewritten as

$$Z_\beta(q) = \sum_{\{\sigma_i\}} \prod_{i<j} (1 + (\exp(\frac{\beta}{N}) - 1) \delta(\sigma_i, \sigma_j)) \exp(\beta h \sum_{\sigma=0}^{q-1} f_\sigma \sum_i \delta(\sigma_i, \sigma)). \quad (3.5.105)$$

The relationship of this expression with graphs appears when we expand the product in this expression. Each of the  $2^{N(N-1)/2}$  terms corresponds to a graph with  $N$  vertices that has an edge between the vertices  $i$  and  $j$  precisely when they both appear in the corresponding term and when  $\sigma_i = \sigma_j$ , that is, when the corresponding Kronecker delta takes the value 1. This graph will in general have several components, but the spin values  $\sigma_i$  are constant on each component by our construction. We can then write the partition function as a sum over graphs,

$$Z_\beta(q) = \sum_{\Gamma} \left( \exp\left(\frac{\beta}{N}\right) - 1 \right)^{l(\Gamma)} \prod_{n=0}^{m(\Gamma)-1} \left( \sum_{\sigma} \exp(\beta h f_{\sigma} S_n) \right) \quad (3.5.106)$$

where  $S_n$  is the size of the  $n$ th component and the product extends over the components of  $\Gamma$ . In order to relate this to our ensemble of random graphs, we put  $\beta = pN$ . Since  $\exp(\frac{\beta}{N}) - 1 = \frac{\beta}{N} + O(\frac{1}{N^2})$ , we can approximate this for large  $N$  as (for  $h = 0$ )

$$Z_{pN}(q) = \sum_{\Gamma} p^{l(\Gamma)} q^{m(\Gamma)} = \exp\left(\frac{pN}{2}\right) z(q) \quad (3.5.107)$$

in leading order in  $N$ , by (3.5.99), (3.5.102). Now, the Potts model exhibits a phase transition to a spontaneous magnetization, that is, all the spins become aligned, above a certain critical value of  $\beta$ . The preceding result then relates this to the appearance of a giant component in our random graph  $\Gamma$  when  $pN$  exceeds a critical threshold. The latter is called a percolation phenomenon, and it is thus related to a phase transition in a statistical mechanics model. The parameter  $h$  in the latter model is useful for deriving properties by taking derivatives at  $h = 0$  whereas the parameter  $q$  becomes useful when one studies large deviation properties, that is the properties of atypical members of our ensemble.

We now generalize the construction of Erdős-Rényi by allowing for different connection probabilities for different pairs of vertices. A generalized random graph is characterized by its number  $N$  of vertices and real numbers  $0 \leq p_{ij} \leq 1$  (with the symmetry  $p_{ij} = p_{ji}$ ) that assign to each pair  $i, j$  of vertices the probability for finding an edge between them. Self-connections of the vertex  $i$  are excluded by  $p_{ii} = 0$ . The expected degree of  $i$  then is

$$\nu_i = \sum_j p_{ij}. \quad (3.5.108)$$

A special case which includes scale-free graphs, is the one of [2]. One starts with an  $N$ -tupel  $\nu = (\nu_1, \dots, \nu_N)$  of positive numbers satisfying

$$\max_i \nu_i^2 \leq \sum_j \nu_j; \quad (3.5.109)$$



when the  $\nu_i$  are positive integers, this is the necessary and sufficient for the existence of a graph with nodes  $i$  of degree  $\nu_i$ ,  $i = 1, \dots, N$ . When putting  $\gamma := \frac{1}{\sum_i \nu_i}$  and  $p_{ij} := \gamma \nu_i \nu_j$ , then  $0 \leq p_{ij} \leq 1$  for all  $i, j$ . We then insert an edge between the nodes  $i$  and  $j$  with probability  $p_{ij}$  to construct the (generalized) random graph  $\Gamma$ . By (3.5.108), the expected degree of node  $i$  in such a graph is  $\nu_i$ . When all the  $\nu_i$  are equal, we obtain an Erdős-Rényi graph. For other types, the degree distribution, i.e., the number of nodes  $i$  with  $\nu_i = k$  will decay as a function of  $k$ , at least for large  $k$ , for example exponentially. When that number behaves like a power  $k^{-\beta}$  instead, we obtain a so-called scale free graph. In the scale-free case, there are thus comparatively more hubs, that is, nodes with large degrees, than in the exponential case.



## Chapter 4

# Pattern formation

We consider spatiotemporal structure formation from interactions between states  $f(x, t)$  at points  $x$  at times  $t$ . This means that the state  $f(x_0, t_0)$  is a function of states  $f(x, t)$  at some or all other points  $x$  at previous times  $t \leq t_0$ , or at least is influenced by some of those states. Here, space, time, and state space can be discrete or continuous. Discreteness or continuity can lead to rather different effects and difficulties. Perhaps that difference is smallest for space. At the appropriate level of abstraction, discrete and continuous space can be treated in the same manner, although some technical aspects are substantially more difficult in the continuous case. In the discrete case, one assumes some underlying graph structure that incorporates which other points  $y$  are the neighbors of a point  $x_0$  whose states  $f(y, t)$  then can directly affect the state  $f(x_0, t_0)$  for  $t_0 \geq t$ . In the continuous case, we need a topology on our space that gives us some notion of infinitesimal proximity for setting up partial differential equations as an analytical framework for pattern formation. Concerning time, the discrete case is usually more difficult than the continuous one. In the latter case, we can work with differential equations whereas in the former one we obtain difference equations or functional iterations. Concerning state space, in the continuous setting we have the possibility of incremental state updates, in particular in the case where time is continuous as well. The discrete case, on the other hand, is more suitable for simulations.

### 4.1 Partial differential equations

Partial differential equations, PDEs for short, constitute a field of mathematics that is distinguished from most other mathematical fields by the fact that a definition of its basic object, a partial differential equation, at best is useless and at worst is severely misleading. In order to understand the essence of this field, one rather needs to study prototypical examples, admitting that what constitutes such a prototype is not clearly defined either. Instead of entering into any further generalities, we start with the perhaps most fundamental one,

the Laplace equation (although this equation is not directly useful in biology). For a twice differentiable function  $u : \Omega \rightarrow \mathbb{R}$  on an open and connected  $\Omega \subset \mathbb{R}^d$ , the Laplacian at  $x = (x^1, \dots, x^d) \in \Omega$  is defined as

$$\Delta u(x) := \sum_{i=1}^d \frac{\partial^2 u}{(\partial x^i)^2}(x). \quad (4.1.1)$$

In the sequel, we shall often abbreviate derivatives by subscripts, i.e.,

$$u_{x^i} := \frac{\partial u}{\partial x^i}, \quad u_{x^i x^i} := \frac{\partial^2 u}{(\partial x^i)^2} \text{ and so on.} \quad (4.1.2)$$

Thus,

$$\Delta u(x) = \sum_{i=1}^d u_{x^i x^i}. \quad (4.1.3)$$

We have already introduced the Laplace operator for a function  $u$  on a graph  $\Gamma$  above, in (2.1.3),

$$\Delta u(x) := \frac{1}{b_x} \left( \sum_{y, y \sim x} u(y) - n_x u(x) \right) \quad (4.1.4)$$

where the vertices  $y$  with  $y \sim x$  are the neighbors of the vertex  $x$ ,  $n_x$  is the degree of  $x$ , that is, the number of its neighbors, and  $b_x$  is a positive factor which we preferred to choose as  $b_x = n_x$ . In order to understand the relationship between these two operators, we replace the domain  $\Omega$  by its discrete approximation by a grid, as is done for example in numerical schemes for solving PDEs. In order not to have to worry about boundary points, we consider for simplicity the case where  $\Omega$  is the entire space  $\mathbb{R}^d$ . For  $h > 0$ , we then define the discrete space

$$\mathbb{R}_h^d := \{(n_1 h, \dots, n_d h)\}, n_1, \dots, n_d \in \mathbb{Z}. \quad (4.1.5)$$

The second partial derivative  $u_{x^i x^i}$  then is approximated by the difference

$$u_{ii} := \frac{1}{h^2} (u(x^1, \dots, x^i + h, \dots, x^d) + u(x^1, \dots, x^i - h, \dots, x^d) - 2u(x^1, \dots, x^i, \dots, x^d)), \quad (4.1.6)$$

and  $\Delta u$  then is approximated by

$$\Delta_h := \sum_{i=1}^d u_{ii}. \quad (4.1.7)$$

When we consider the grid  $\mathbb{R}_h^d$  as a graph on which the neighbors of  $x = (x^1, \dots, x^d)$  are the points  $(x^1, \dots, x^i \pm h, \dots, x^d)$ ,  $i = 1, \dots, d$ , up to a factor, this is the same as the graph Laplacian (the factor  $\frac{1}{h^2}$  has been chosen here in order that the discrete Laplacian converges to the continuous one for  $h \rightarrow 0$ ).

**Definition 4.1.1.** A function  $u$  (on a domain  $\Omega$  or a graph  $\Gamma$ ) is called harmonic if it satisfies the Laplace equation

$$\Delta u = 0. \quad (4.1.8)$$

In the discrete case, from (4.1.4) it is clear that a harmonic function  $u$  satisfies the mean value property

$$u(x) = \frac{1}{n_x} \sum_{y, y \sim x} u(y) \quad (4.1.9)$$

for all  $x$ . The mean value property of harmonic functions also holds in the continuous case:

$$u(x) = \frac{1}{\omega_d r^d} \int_{B(x,r)} u(y) dy \quad (4.1.10)$$

whenever the ball  $B(x, r) := \{y \in \mathbb{R}^d : \|y - x\| < r\}$  of radius  $r$  around  $x$  is contained in  $\Omega$ . Here,  $\omega_d$  is the volume of the unit ball in  $\mathbb{R}^d$ . Conversely, one can show that the mean value property for a continuous  $u$  implies that it is harmonic. More generally,  $u$  is called subharmonic in  $\Omega$  if

$$\Delta u(x) \geq 0 \text{ for } x \in \Omega. \quad (4.1.11)$$

This turns out to be equivalent to the mean value inequality

$$u(x) \leq \frac{1}{\omega_d r^d} \int_{B(x,r)} u(y) dy. \quad (4.1.12)$$

From the mean value property, one easily derives the maximum principle:

**Lemma 4.1.1.** *Suppose that  $u$  is harmonic, or more generally, subharmonic in the open and connected  $\Omega$ . If there exists some  $x_0 \in \Omega$  with*

$$u(x_0) = \sup_{x \in \Omega} u(x), \quad (4.1.13)$$

*then  $u$  is constant in  $\Omega$ . This is the so-called strong maximum principle, and it implies the weak maximum principle: If  $\Omega$  is bounded and  $u \in C^0(\bar{\Omega})$  (meaning that  $u$  is defined and continuous on the closure of  $\Omega$ ), then*

$$u(x) \leq \max_{y \in \partial\Omega} u(y). \quad (4.1.14)$$

*Finally, if a nonconstant  $u$  assumes its maximum at the smooth boundary point  $y_0$  and if it is differentiable there, then*

$$\frac{\partial}{\partial n} u(y_0) > 0 \quad (4.1.15)$$

*where  $\frac{\partial}{\partial n}$  denotes the derivative in the direction of the exterior normal of  $\Omega$ .*

*Proof.* When  $u(x_0) = \sup_{x \in \Omega} u(x) =: M$ , we put

$$\Omega^M := \{y \in \Omega : u(y) = M\} \neq \emptyset.$$

For  $z \in \Omega^M$  with  $r > 0$  such that  $B(z, r) \subset \Omega$ , we get

$$0 = u(z) - M \leq \frac{1}{\omega_d r^d} \int_{B(z,r)} (u(y) - M) dy \leq 0 \quad (4.1.16)$$

since  $M$  is the supremum of  $u$ , and we see that necessarily  $u(y) = M$  for all  $y \in B(z, r)$ . Therefore, whenever  $z \in \Omega^M$  and  $B(z, r) \subset \Omega$ , then that entire ball is also contained in  $\Omega^M$ . Since  $\Omega$  is connected,  $\Omega^M$  has to be all of  $\Omega$ . This means that  $u \equiv M$  in  $\Omega$  which is what we wanted to prove. The weak maximum principle then follows from the simple observation that a continuous function on the bounded and closed, hence compact set  $\bar{\Omega}$  has to assume its supremum. When a harmonic function does so in the interior  $\Omega$ , it is constant by the strong maximum principle, and (4.1.14) holds, and when the supremum is assumed on the boundary, (4.1.14) holds as well. For the proof of the boundary point maximum result, we refer to the literature, e.g. [21].  $\square$

The weak maximum principle can also be expressed by saying that a non-constant harmonic function assumes its supremum only on the boundary of  $\Omega$  when that set is bounded and  $u$  is continuous on the closure of  $\Omega$ .

The strong maximum principle also holds in the discrete case, with the same kind of proof. Since a graph is an object without a boundary, this implies that any harmonic function on a finite graph is constant. Of course, one can also turn the situation into a boundary value problem by declaring a subset  $S_0$  of the vertex set  $S$  of  $\Gamma$  as the boundary and considering the problem

$$\Delta u(x) = 0 \text{ for } x \in S \setminus S_0 \quad (4.1.17)$$

$$u(x) = g(x) \text{ for } x \in S_0 \quad (4.1.18)$$

for some prescribed function  $g : S_0 \rightarrow \mathbb{R}$ .

By the mean value formula, harmonic functions represent equilibrium states where the value at each point is the average of the values of its neighbors. This observation also suggests a scheme for the proof of the existence of harmonic functions, for example for given boundary values  $g$  on  $\partial\Omega$ . One starts with any (continuous) function  $u_0 : \bar{\Omega} \rightarrow \mathbb{R}$  with  $u_0 = g$  on  $\partial\Omega$ . Having constructed  $u_1, \dots, u_{n-1}$  iteratively, one finds  $u_n(x)$  for  $x \in \Omega$  by replacing  $u_{n-1}(x)$  by its mean value on some ball  $B(x, r) \subset \Omega$ . This simple idea can be made to work, and it yields a constructive scheme for finding a harmonic  $u$  with given boundary values. On a graph  $\Gamma$ , this means

$$u(x, t+1) := \frac{1}{n_x} \sum_{y \sim x} u(y, t) \quad (4.1.19)$$

for  $x \in S - S_0, t \in \mathbb{N}$ ,  $u(x, 0)$  being an arbitrary function satisfying the boundary condition (4.1.17). For the numerical implementation (for using this for solving the boundary value problem for harmonic functions in a continuous domain by discrete approximation), one again replaces  $\Omega$  by a discrete grid of some small mesh size  $h$ . For temporal step size  $k$ , one puts

$$u(x; t+k) := \frac{1}{2d} \sum_{i=1}^d (u(x^1, \dots, x^i - h, \dots, x^d; t) + u(x^1, \dots, x^i + h, \dots, x^d; t)) \quad (4.1.20)$$

Conceptually, this is quite useful because it suggests a PDE that is defined on space and time instead of on space only as the Laplace equation and that models the approach to equilibrium. This is the heat equation

$$\frac{\partial}{\partial t}u(x, t) = \Delta u(x, t) \quad (= \sum_{i=1}^d \frac{\partial^2}{(\partial x^i)^2}u(x, t)) \quad (4.1.21)$$

or abbreviated,

$$u_t(x, t) = \Delta u(x, t). \quad (4.1.22)$$

In fact, the straightforward discretization of (4.1.22) is

$$\begin{aligned} & \frac{1}{k}(u(x, t+k) - u(x, t)) & (4.1.23) \\ = & \frac{1}{h^2} \sum_{i=1}^d (u(x^1, \dots, x^i - h, \dots, x^d, t) + u(x^1, \dots, x^i + h, \dots, x^d, t) - 2u(x^1, \dots, x^d, t)). \end{aligned}$$

For  $2dk = h^2$ , the term  $u(x, t)$  cancels in (4.1.23), and we obtain (4.1.20).<sup>1</sup>

For the heat equation, we also have a maximum principle:

**Lemma 4.1.2.** *Let  $\Omega \subset \mathbb{R}^d$  be open,  $0 < T \leq \infty$ , and let  $u(x, t)$  be continuous for  $x \in \bar{\Omega}, 0 \leq t \leq T$  and satisfy*

$$u_t(x, t) = \Delta u(x, t) \text{ for } x \in \Omega, 0 < t < T. \quad (4.1.24)$$

Then

$$\sup_{\bar{\Omega} \times [0, T]} u = \sup_{(\bar{\Omega} \times \{0\}) \cup (\partial\Omega \times [0, T])} u. \quad (4.1.25)$$

When  $T < \infty$ , the supremum becomes a maximum. Again, there is also a strong version of the maximum principle, saying that a solution of (4.1.24) cannot attain a maximum in  $\Omega \times (0, T]$  without being constant. Also, there is an analogue of the result for boundary maxima, that is, at a nontrivial boundary point maximum, one obtains a positive exterior normal derivative, if the situation is sufficiently smooth.

The Lemma says that the maximum of a solution of the heat equation is always attained either at the spatial boundary  $\partial\Omega \times [0, T]$  of the cylinder  $\Omega \times (0, T)$  or at the initial set  $t = 0$ .

In the continuous as in the discrete case, one shows that a solution of the initial boundary value for the heat equation

$$u_t(x, t) = \Delta u(x, t) \text{ for } x \in \Omega, 0 < t \quad (4.1.26)$$

$$u(x, 0) = u_0(x) \text{ for } x \in \Omega \quad (4.1.27)$$

$$u(y, t) = g(y) \text{ for } y \in \partial\Omega \quad (4.1.28)$$

---

<sup>1</sup>The fact that the temporal step size  $k$  satisfies  $2dk = h^2$  slows down the convergence of the scheme for  $h \rightarrow 0$  and makes this not really a good numerical scheme in practice.

for continuous initial values  $u_0$  and boundary values  $g$  (and under certain mild technical assumptions on  $\Omega$ ) converges to a solution of the boundary value problem for the Laplace equation, the Dirichlet problem

$$\Delta u(x) = 0 \text{ for } x \in \Omega \quad (4.1.29)$$

$$u(y) = g(y) \text{ for } y \in \partial\Omega \quad (4.1.30)$$

for  $t \rightarrow \infty$ , that is,  $\lim_{t \rightarrow \infty} u(x, t) = u(x)$ .

A generalization of the Laplace equation is the Poisson equation

$$\Delta u(x) = f(x) \text{ for } x \in \Omega \quad (4.1.31)$$

for some given function  $f : \Omega \rightarrow \mathbb{R}$ , or its analogues in the discrete case or for the heat equation.

A basic idea for solving (4.1.31) consists in the superposition of point solutions. That means that for each  $y \in \Omega$ , we try to find some function  $\gamma(x, y)$  that solves (4.1.31) at  $y$  and is harmonic elsewhere. If we then integrate w.r.t.  $y$ , we should obtain the desired solution of (4.1.31). Let us first try to implement this in the discrete case where we only have to take a sum in place of an integral. Of course, instead of  $f(y)$ , we can then take 1 as the right hand side for our equation at  $y$  and multiply the result by  $f(y)$  and then sum w.r.t.  $y$ . Thus, given a graph  $\Gamma$  as before, and a node  $y \in S$ , the vertex set of  $\Gamma$ , we want to solve

$$\Delta_x G(x, y) (= \frac{1}{n_x} \sum_{z, z \sim x} (G(z, y) - G(x, y))) = \frac{1}{n_x} \delta(x, y) \quad (4.1.32)$$

with

$$\delta(x, y) := \begin{cases} 1 & \text{for } x = y \\ 0 & \text{elsewhere.} \end{cases} \quad (4.1.33)$$

(Since we are looking at functions of two variables, we indicate by a subscript w.r.t. which variable the Laplacian  $\Delta$  acts.) A solution of (4.1.32) is called a Green function.

If we can find such a Green function, a solution to the discrete Poisson equation

$$\Delta u(x) = f(x) \text{ for } x \in S \quad (4.1.34)$$

is then simply given by

$$u(x) = n_x \sum_y G(x, y) f(y). \quad (4.1.35)$$

There is one problem here: We cannot solve (4.1.32) because for any function  $g$  on a graph  $\Gamma$ , we have

$$\sum_x n_x \Delta g(x) = 0 \quad (4.1.36)$$

and therefore necessarily also

$$\sum_x n_x \Delta_x G(x, y) = 0, \quad (4.1.37)$$



but the right hand side of (4.1.32) does not fulfill that condition. In abstract terms, the Laplacian has a kernel, consisting of the constant functions, and is therefore not invertible. It is invertible only on the space orthogonal to the constants. That means that we can expect to solve (4.1.34) only when  $f$  satisfies

$$\sum_x n_x f(x) = 0. \quad (4.1.38)$$

This can be easily remedied, however. We simply replace (4.1.32) by

$$\Delta_x G(x, y) = \frac{1}{n_x} \delta(x, y) - \frac{1}{\sum_z n_z}, \quad (4.1.39)$$

that is, subtract a suitable constant on the right hand side so as to achieve (4.1.37). When (4.1.38) holds, the contribution of the constant disappears in (4.1.35), and so, we can solve (4.1.34) for those  $f$ .

Another possibility to circumvent that problem is to impose a boundary condition, that is, solve

$$\Delta u(x) = f(x) \text{ for } x \in S \setminus S_0 \quad (4.1.40)$$

$$u(x) = g(x) \text{ for } x \in S_0 \quad (4.1.41)$$

for some prescribed function  $g : S_0 \rightarrow \mathbb{R}$ . Here, we assume  $S_0 \neq \emptyset$ , but otherwise,  $S_0$  is completely arbitrary.

In order to achieve that, we first consider the homogeneous boundary condition, that is,  $g = 0$ . For that, we impose the homogeneous boundary condition

$$G(x, y) = 0 \text{ for } x \in S_0 \text{ and all } y \quad (4.1.42)$$

take the corresponding  $u$  from (4.1.35) (the equation now imposed for  $x \in S \setminus S_0$ ), which then satisfies  $u(x) = 0$  for  $x \in S_0$ . In order to solve the general boundary value problem, we then simply add a solution  $u_0(x)$  of (4.1.17) to get the right boundary condition. In abstract terms, imposing a boundary condition eliminates the kernel of the Laplacian. We can then not only solve the boundary value problem, but the solution is also unique, because the difference of two solutions is a harmonic function with zero boundary values, hence identically zero itself (as follows in many ways, for example from the maximum principle). In the continuous case, we can use the same strategy. We want to solve

$$\Delta u(x) = f(x) \text{ for } x \in \Omega \quad (4.1.43)$$

$$u(x) = g(x) \text{ for } x \in \partial\Omega. \quad (4.1.44)$$

Again, assuming that we can already solve the boundary value problem for the Laplace equation, that is, find a solution for

$$\Delta u(x) = 0 \text{ for } x \in \Omega \quad (4.1.45)$$

$$u(x) = g(x) \text{ for } x \in \partial\Omega, \quad (4.1.46)$$

we consider homogenous boundary values, that is,

$$\Delta u(x) = f(x) \text{ for } x \in \Omega \quad (4.1.47)$$

$$u(x) = 0 \text{ for } x \in \partial\Omega. \quad (4.1.48)$$

As before, we start with

$$\Delta_x G(x, y) = \delta(x, y), \quad (4.1.49)$$

the Dirac delta functional. This means that for every continuous  $\phi$ , we have

$$\phi(x) = \int_{\Omega} \delta(x, y)\phi(y)dy = \int_{\Omega} \Delta_x G(x, y)\phi(y)dy. \quad (4.1.50)$$

And

$$u(x) = \int_{\Omega} G(x, y)f(y)dy \quad (4.1.51)$$

then satisfies

$$\Delta_x u(x) = \int_{\Omega} \Delta_x G(x, y)f(y)dy = f(x). \quad (4.1.52)$$

Once more, in order to get homogeneous boundary values, that is,  $u|_{\partial\Omega} = 0$ , for  $u$  in (4.1.51), we need to have  $G(x, y) = 0$  for  $x \in \partial\Omega$ . This can indeed be achieved, but we do not go into the details here. We rather display the so-called fundamental solutions, particular solutions of (4.1.49) in the whole space  $\mathbb{R}^d$ . These are

$$\Gamma(x, y) = \Gamma(|x - y|) := \begin{cases} \frac{1}{2\pi} \log |x - y| & \text{for } d = 2 \\ \frac{1}{d(2-d)\omega_d} |x - y|^{2-d} & \text{for } d > 2 \end{cases} \quad (4.1.53)$$

where  $\omega_d$  is the volume of the  $d$ -dimensional unit ball  $B(0, 1) \subset \mathbb{R}^d$ . The computations that this  $\Gamma$  solves (4.1.49) are straightforward, but somewhat lengthy and omitted here.

For the heat equation, we also have a fundamental solution from which more general problems can be solved by superposition. For  $x, y \in \mathbb{R}^d, t > 0$ , we put

$$K(x, y, t) := \frac{1}{(4\pi t)^{d/2}} e^{-\frac{|x-y|^2}{4t}}. \quad (4.1.54)$$

$K$  solves the heat equation:

$$\frac{\partial}{\partial t} K(x, y, t) = \Delta K(x, y, t) \text{ for all } x, y \in \mathbb{R}^d, t > 0. \quad (4.1.55)$$

We have the normalization

$$\int_{\mathbb{R}^d} K(x, y, t)dy = 1 \text{ for all } x \in \mathbb{R}^d, t > 0. \quad (4.1.56)$$

Also, for a bounded and continuous function  $f$  on  $\mathbb{R}^d$ ,

$$u(x, t) = \int_{\mathbb{R}^d} K(x, y, t)f(y)dy \quad (4.1.57)$$

solves the heat equation

$$u_t = \Delta u \quad (4.1.58)$$

for  $x \in \mathbb{R}^d, t > 0$  and has the initial values

$$\lim_{t \rightarrow 0} u(x, t) = f(x) \quad (4.1.59)$$

which is abbreviated as

$$u(x, 0) = f(x). \quad (4.1.60)$$

We now briefly discuss the eigenvalue problem for the Laplace operator and its connections with the heat equation. Again, this is formally analogous to the discrete case, already treated in 2.1.3, although the details now require a more careful analysis and depend on some analytical result, the Rellich compactness theorem.

The eigenvalue problem for the Laplace operator consists in finding nontrivial solutions of

$$\Delta u(x) + \lambda u(x) = 0 \quad \text{in } \Omega, \quad (4.1.61)$$

for some constant  $\lambda$ , the eigenvalue in question. Here one also imposes some boundary conditions on  $u$ . It seems natural to require the Dirichlet boundary condition

$$u = 0 \quad \text{on } \partial\Omega. \quad (4.1.62)$$

For many applications, however, it is more natural to have the Neumann boundary condition

$$\frac{\partial u}{\partial n} = 0 \quad \text{on } \partial\Omega \quad (4.1.63)$$

instead, where  $\frac{\partial}{\partial n}$  denotes the derivative in the direction of the exterior normal. Here, in order to make this meaningful, one needs to impose suitable regularity of  $\partial\Omega$ . For simplicity, we shall assume that  $\Omega$  is a  $C^\infty$ -domain in treating Neumann boundary conditions. When the domain is a closed manifold in place of a subset with boundary of  $\mathbb{R}^d$ , one does not impose any further condition, as in the case of a graph.

We shall employ the  $L^2$ -product

$$\langle f, g \rangle := \int_{\Omega} f(x)g(x)dx \quad (4.1.64)$$

for  $f, g \in L^2(\Omega)$ , that is,  $\int_{\Omega} f(x)^2 dx, \int_{\Omega} g(x)^2 dx < \infty$  and we shall also put

$$\|f\| := \|f\|_{L^2(\Omega)} = \langle f, f \rangle^{\frac{1}{2}}. \quad (4.1.65)$$

We note the symmetry of the Laplace operator,

$$\langle \Delta\varphi, \psi \rangle = -\langle D\varphi, D\psi \rangle = \langle \varphi, \Delta\psi \rangle \quad (4.1.66)$$

for all  $\varphi, \psi \in C_0^\infty(\Omega)$ , as well as for  $\varphi, \psi \in C^\infty(\Omega)$  with  $\frac{\partial \varphi}{\partial n} = 0 = \frac{\partial \psi}{\partial n}$  on  $\partial\Omega$ . Here,  $D\varphi$  abbreviates the vector  $\frac{\partial}{\partial x^1}\varphi, \dots, \frac{\partial}{\partial x^d}\varphi$ . This symmetry implies that all eigenvalues are real.

**Theorem 4.1.1.** *Let  $\Omega \subset \mathbb{R}^d$  be open and bounded. Then the eigenvalue problem*

$$\Delta u + \lambda u = 0, \quad u = 0 \text{ on } \partial\Omega$$

*has countably many eigenvalues*

$$0 < \lambda_1 < \lambda_2 \leq \dots \leq \lambda_m \leq \dots \quad (4.1.67)$$

*with*

$$\lim_{m \rightarrow \infty} \lambda_m = \infty$$

*and pairwise  $L^2$ -orthonormal eigenfunctions  $u_i$  and  $\langle Du_i, Du_i \rangle = \lambda_i$ . Any  $v \in L^2(\Omega)$  can be expanded in terms of these eigenfunctions,*

$$v = \sum_{i=1}^{\infty} \langle v, u_i \rangle u_i \quad (\text{and thus } \langle v, v \rangle = \sum_{i=1}^{\infty} \langle v, u_i \rangle^2). \quad (4.1.68)$$

*Moreover, the first eigenfunction  $u_1$  does not change sign in  $\Omega$ , that is, we may assume*

$$u_1 > 0 \text{ in } \Omega. \quad (4.1.69)$$

We should explain the inequality signs in (4.1.67). An eigenvalue can have higher multiplicity, meaning that there may exist several linearly independent eigenfunctions with the same eigenvalue. Therefore, eigenvalues are counted according to the dimension of the eigenspaces. The inequality  $0 < \lambda_1$  simply says that the first eigenvalue is positive. The inequality  $\lambda_1 < \lambda_2$  represents the theorem that the first eigenvalue is simple, that is, the eigenspace corresponding to  $\lambda_1$  is one-dimensional.

For Neumann boundary conditions, we have an analogous result:

**Theorem 4.1.2.** *Let  $\Omega \subset \mathbb{R}^d$  be bounded, open, and of class  $C^\infty$ . Then the eigenvalue problem*

$$\Delta u + \lambda u = 0$$

*has countably many eigenvalues*

$$0 = \lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_m \leq \dots$$

*with*

$$\lim_{n \rightarrow \infty} \lambda_n = \infty$$

and pairwise  $L^2$ -orthonormal eigenfunctions  $u_i$  that satisfy

$$\frac{\partial u_i}{\partial n} = 0 \quad \text{on } \partial\Omega.$$

Any  $v \in L^2(\Omega)$  can be expanded in terms of these eigenfunctions

$$v = \sum_{i=0}^{\infty} \langle v, u_i \rangle u_i \quad (\text{and thus } \langle v, v \rangle = \sum_{i=0}^{\infty} \langle v, u_i \rangle^2). \quad (4.1.70)$$

In Theorem 4.1.2,  $\lambda_0 = 0$  appears as an eigenvalue. In fact, any non-vanishing constant is an eigenfunction with eigenvalue 0, and, in contrast to the Dirichlet condition, these are not excluded by the Neumann boundary condition. When  $\Omega$  has more than one component, we can in fact choose a different constant on each component. When  $\Omega$  is connected, however, a global constant is the only eigenfunction with eigenvalue 0, and this then is a simple eigenvalue. It is also insightful and instructive to look at the scaling behavior of the eigenvalues. If instead of  $\Omega$  we consider the domain  $\alpha\Omega := \{\alpha x : x \in \Omega\}$  for a scaling factor  $\alpha > 0$ , then its eigenfunctions are given by  $u_i(\frac{y}{\alpha})$ . Since  $\frac{\partial^2}{(\partial y^i)^2} u(\frac{y}{\alpha}) = \frac{\partial^2}{(\partial x^i)^2} u(x)$  for  $y^i = \alpha x^i$ , the eigenvalues of  $\Omega_\alpha$  then are  $\alpha^{-2} \lambda_i$  where the  $\lambda_i$ , of course, are the ones of  $\Omega$  (this argument is valid for both Dirichlet and Neumann eigenvalues). Since the volume  $\|\Omega_\alpha\|$  of  $\Omega_\alpha$  is  $\alpha^d \|\Omega\|$  for a  $d$ -dimensional domain, the eigenvalues scale like  $\text{Vol}^{-\frac{2}{d}}$ . The Weyl type estimates state that (under some mild regularity assumptions on  $\Omega$ .) the eigenvalues  $\lambda_k$  of  $\Omega$  grow proportionally to  $(\frac{k}{\|\Omega\|})^{\frac{2}{d}}$  up to terms of lower order.

Below, we shall also need the following result. We consider the average of  $v$  on  $\Omega$

$$\bar{v} := \frac{1}{\|\Omega\|} \int_{\Omega} v(x) dx \quad (4.1.71)$$

where  $\|\Omega\|$  is the volume of  $\Omega$ .

**Corollary 4.1.1.** *Let  $\lambda_1$  be the first Neumann eigenvalue of  $\Omega$ .*

*For  $v \in H^{1,2}(\Omega)$  (that is, it not only is square integrable itself, but also has square integrable first derivatives in the  $L^2$ -sense) with  $\frac{\partial v}{\partial n}$  on  $\partial\Omega$*

$$\lambda_1 \langle v - \bar{v}, v - \bar{v} \rangle \leq \langle Dv, Dv \rangle. \quad (4.1.72)$$

*For  $v$  of class  $H^{2,2}(\Omega)$  (that is, it also has square integrable second derivatives in the  $L^2$ -sense), also*

$$\lambda_1 \langle Dv, Dv \rangle \leq \langle \Delta v, \Delta v \rangle. \quad (4.1.73)$$

*Proof.* We observe that  $Dv = D(v - \bar{v})$ ,  $\Delta v = \Delta(v - \bar{v})$ , and

$$\langle v - \bar{v}, v - \bar{v} \rangle = \sum_{i=1}^{\infty} \langle v, u_i \rangle^2, \quad (4.1.74)$$

that is, the term for  $i = 0$  disappears from the expansion because  $v - \bar{v}$  is orthogonal to the constant eigenfunction  $u_0$ . Using

$$\begin{aligned}\langle Dv, Dv \rangle &= \sum_{i=1}^{\infty} \lambda_i \langle v, u_i \rangle^2 \\ \langle \Delta v, \Delta v \rangle &= \sum_{i=1}^{\infty} \lambda_i^2 \langle v, u_i \rangle^2\end{aligned}$$

and  $\lambda_1 \leq \lambda_i$  then yields (4.1.72), (4.1.73). □ □

**Remark:** The following argument that assumes still more regularity of  $v$  is also instructive:

$$\int_{\Omega} (\Delta v)^2 = \int \sum_{i=1}^d v_{x^i x^i} \sum_{j=1}^d v_{x^j x^j} = \int \sum_{i,j=1}^d v_{x^i x^j} v_{x^i x^j}, \quad (4.1.75)$$

integrating by parts twice, without incurring a boundary term because of the Neumann boundary condition. Therefore, applying (4.1.72) to  $v_{x^i}$  for  $i = 1, \dots, d$  yields (4.1.73).

One can use the eigenfunctions of the Laplacian to write an expansion for the Green function. We consider the case of the Dirichlet boundary conditions as in theorem 4.1.1. Thus, the Green function has to solve

$$\Delta_x G(x, y) = \delta(x, y) \text{ for } x, y \in \Omega \quad (4.1.76)$$

$$G(x, y) = 0 \text{ for } x \in \partial\Omega, y \in \Omega. \quad (4.1.77)$$

This Green function can then be represented in terms of the Dirichlet eigenfunctions of theorem 4.1.1 as

$$G(x, y) = - \sum_n \frac{1}{\lambda_n} u_n(x) u_n(y). \quad (4.1.78)$$

To see this, recalling (4.1.51), (4.1.52), we consider

$$u(x) = - \int_{\Omega} \sum_n \frac{1}{\lambda_n} u_n(x) u_n(y) f(y) dy \quad (4.1.79)$$

and compute

$$\begin{aligned}\Delta_x u(x) &= \sum_n u_n(x) \int u_n(y) f(y) dy \\ &= \sum_n u_n(x) \langle u_n, f \rangle \\ &= f(x) \text{ by (4.1.68).}\end{aligned}$$

These expansions in terms of eigenfunctions of the Laplace operator are also useful for the heat equation

$$u_t(x, t) = \Delta u(x, t) \quad \text{for } x \in \Omega, 0 < t. \quad (4.1.80)$$

We try to find solutions with separated variables, i.e., of the form

$$u(x, t) = v(x)w(t). \quad (4.1.81)$$

Inserting this ansatz into (4.1.80), we obtain

$$\frac{w_t(t)}{w(t)} = \frac{\Delta v(x)}{v(x)}. \quad (4.1.82)$$

Since the left-hand side of (4.1.82) is a function of  $t$  only, while the right-hand side is a function of  $x$ , each of them has to be constant. Thus

$$\Delta v(x) = -\lambda v(x), \quad (4.1.83)$$

$$w_t(t) = -\lambda w(t), \quad (4.1.84)$$

for some constant  $\lambda$ . We consider the case where we assume homogeneous boundary conditions on  $\partial\Omega \times [0, \infty)$ , i.e.,

$$u(x, t) = 0 \quad \text{for } x \in \partial\Omega \quad (4.1.85)$$

or equivalently,

$$v(x) = 0 \quad \text{for } x \in \partial\Omega. \quad (4.1.86)$$

A nontrivial solution  $v$  of (4.1.83), (4.1.85) is an eigenfunction of the Laplace operator, and  $\lambda$  an eigenvalue. By theorem 4.1.1, the eigenvalues constitute a discrete sequence  $(\lambda_n)_{n \in \mathbb{N}}$ ,  $\lambda_n \rightarrow \infty$  for  $n \rightarrow \infty$ . Thus, a nontrivial solution of (4.1.83), (4.1.86) exists precisely if  $\lambda = \lambda_n$ , for some  $n \in \mathbb{N}$ . The solution of (4.1.84) then is simply given by

$$w(t) = w(0)e^{-\lambda t}.$$

So, if we denote an eigenfunction for the eigenvalue  $\lambda_n$  by  $u_n$ , we obtain the solution

$$u(x, t) = u_n(x)w(0)e^{-\lambda_n t}$$

of the heat equation (4.1.80), with the homogeneous boundary condition

$$u(x, t) = 0 \quad \text{for } x \in \partial\Omega$$

and the initial condition

$$u(x, 0) = u_n(x)w(0).$$

This seems to be a rather special solution. Nevertheless, in a certain sense this is the prototype of a solution. Namely, because (4.1.80) is a linear equation, any linear combination of solutions is a solution itself, and so we may take sums of such solutions for different eigenvalues  $\lambda_n$ . In fact, by theorem 4.1.1, any  $L^2$ -function on  $\Omega$ , and thus in particular any continuous function  $f$  on  $\bar{\Omega}$ , assuming  $\Omega$  to be bounded, that vanishes on  $\partial\Omega$ , can be expanded as

$$f(x) = \sum_{n \in \mathbb{N}} \alpha_n u_n(x), \quad (4.1.87)$$

where the  $u_n(x)$  are the orthonormal eigenfunctions of  $\Delta$ ,

$$\int_{\Omega} u_n(x)u_m(x)dx = \delta_{nm},$$

and with

$$\alpha_n = \int_{\Omega} u_n(x)f(x)dx.$$

We then have an expansion for the solution of

$$\begin{aligned} u_t(x, t) &= \Delta u(x, t) && \text{for } x \in \Omega, t \geq 0, \\ u(x, t) &= 0 && \text{for } x \in \partial\Omega, t \geq 0, \\ u(x, 0) &= f(x) && \left( = \sum_n \alpha_n u_n(x) \right), \quad \text{for } x \in \Omega, \end{aligned} \quad (4.1.88)$$

namely,

$$u(x, t) = \sum_{n \in \mathbb{N}} \alpha_n e^{-\lambda_n t} u_n(x). \quad (4.1.89)$$

Since all the  $\lambda_n$  are nonnegative, we see from this representation that all the “modes”  $\alpha_n u_n(x)$  of the initial values  $f$  are decaying in time for a solution of the heat equation. In this sense, the heat equation regularizes or smoothes out its initial values. In particular, since thus all factors  $e^{-\lambda_n t}$  are less than or equal to 1 for  $t \geq 0$ , the series (4.1.89) converges in  $L^2(\Omega)$ , because (4.1.87) does.

If we write

$$q(x, y, t) := \sum_{n \in \mathbb{N}} e^{-\lambda_n t} u_n(x)u_n(y), \quad (4.1.90)$$

theorem 4.1.1 shows convergence of this series, and we may represent the solution  $u(x, t)$  of (4.1.88) as

$$\begin{aligned} u(x, t) &= \sum_{n \in \mathbb{N}} e^{-\lambda_n t} u_n(x) \int_{\Omega} u_n(y)f(y)dy && \text{by (4.1.89)} \\ &= \int_{\Omega} q(x, y, t)f(y)dy. \end{aligned} \quad (4.1.91)$$

Comparing this with (4.1.57), (4.1.58), we see that  $q(x, y, t)$  as in (4.1.90) yields a heat kernel, in analogy to formula (4.1.79) for the Laplace equation.

## 4.2 Diffusion and random walks

In this section, we want to explore the relationship between partial differential equations and stochastic analysis. As before, we start with the discrete case. We can partly follow the classical treatment of [7]. We take some graph  $\Gamma$  with



vertex set  $V$ , and some subset  $V_0$  of  $V$  that is considered as the boundary.  $V \setminus V_0$  is the interior. The choice of  $V_0$  is rather arbitrary. It should not be empty, at least on certain occasions, nor should it coincide with the entire vertex set. Although we shall not indicate those places below where such assumptions are used, we might also wish to require that the graph obtained from  $\Gamma$  by eliminating the vertex set  $V_0$  and all edges connected to vertices in  $V_0$  is still connected.

We now construct a diffusion process on  $\Gamma$ . We assume that we have a unit of some substance at the point  $x \in V$ . That substance is diffusing in  $\Gamma$  in such manner that the fraction  $s$  of our substance present at the vertex  $y$  at time  $n$  is equally distributed among the neighbors of  $y$  at time  $n + 1$ , that is, each neighbor of  $y$  receives  $\frac{s}{n_y}$  of the substance where  $n_y$ , as always, is the degree of  $y$ . Whatever amount of the substance reaches a boundary point will stay there forever. When the initial point  $x$  was a boundary point, the whole amount of the substance will stay there. Thus, the boundary is absorbing for our diffusion process. – According to these rules, the total amount of our substance present in  $\Gamma$  at any time  $n$  is always the same, that is, our diffusion process satisfies a conservation law.

There is an alternative view of this process, and identifying those two views will be very insightful. That latter view is the one of a random walk on  $\Gamma$ . This is a stochastic process, discrete Brownian motion, with discrete time  $n \in \mathbb{N}$ . A walker or a Brownian particle, whatever physical interpretation one prefers, starts at  $x$ , and when it happens to be at the interior point  $y$  at time  $n$ , it moves to one of the neighbors of  $y$  at time  $n + 1$ , and all these neighbors have the same probability, that is  $\frac{1}{n_y}$ , of receiving the particle.

We see that the probabilities follow the above diffusion process. Let  $z$  be a boundary point. The probability  $w(x, z)$  of reaching  $z$  by a random walk starting at  $x$  without having previously hit any other boundary point then equals the fraction of our diffusing substance that has accumulated at  $z$  in infinite time. According to our rules at the boundary, we also have

$$w(z, z) = 1 \text{ for } z \in V_0 \quad \text{and } w(z, z') = 0 \text{ when } z \neq z' \in V_0. \quad (4.2.92)$$

When  $z$  is the first boundary point reached by a random walk starting at  $x$ , we also call  $z$  the exit point of that random walk. When  $r_n(x, z)$  is the fraction of our substance reaching  $z$  from  $x$  after precisely  $n$  steps, we have

$$w(x, z) = \sum_{\nu=0}^{\infty} r_{\nu}(x, z). \quad (4.2.93)$$

Likewise, the probability  $q_n(x, y)$  of reaching  $y$  after  $n$  steps starting from  $x$  equals the fraction of our substance that happens to be at time  $n$  at  $y$ , in case  $x$  and  $y$  are interior points. When either of them is a boundary point, that probability is put to 0. We are also interested in the sum

$$v(x, y) := \sum_{\nu=0}^{\infty} q_{\nu}(x, y). \quad (4.2.94)$$

When  $x$  and  $y$  are interior points, this equals the amount of substance that has passed through  $y$  at some time. In the probabilistic interpretation, this is the expected number of times the random walk starting at  $x$  passes through  $y$  before exiting at some boundary point.

The sum in (4.2.93) converges because its members are nonnegative and its partial sums cannot become larger than 1 as only some fraction of the original substance can reach  $z$  before being absorbed at some other boundary point. This then also implies that  $q_n(x, y)$  tends to 0 for  $n \rightarrow \infty$ . Indeed, let  $q_n(x, y) > \epsilon$ , and assume that the boundary point  $z$  can be reached from  $y$  along some path  $y_0 := y, y_1, \dots, y_m = z$  that does not hit the boundary before  $z$ . Then after  $m$  steps, a fraction  $\frac{\epsilon}{n_{y_0} n_{y_1} \dots n_{y_{m-1}}}$  of our substance reaches  $z$  along that particular path and is absorbed at  $z$ . By convergence of the series in (4.2.93), the fraction of substance reaching  $z$  after  $n$  steps tends to 0 for  $n \rightarrow \infty$ . Therefore,  $q_n(x, y)$  also has to converge to 0 for  $n \rightarrow \infty$ . In particular, the probability of staying in the interior for an infinite amount of time vanishes. Consequently,

$$\sum_{z \in V_0} w(x, z) = 1. \quad (4.2.95)$$

We may then consider  $w(x, \cdot)$  as a probability distribution for the exit point of the random walk starting at  $x$ .

We now turn to proving the convergence of the series in (4.2.94). We have the relation

$$q_n(x, y) = \sum_{y' \sim y} \frac{1}{n_{y'}} q_{n-1}(x, y') \quad (4.2.96)$$

for  $n > 1$  because whatever is reaching  $y$  at some time  $n$  has to be at some neighbor  $y'$  of  $y$  at time  $n - 1$ . We recall here that  $q_{n-1}(x, y') = 0$  when  $y'$  happens to be a boundary point. Also,  $q_0(x, x) = 1$  when  $x$  itself is an interior point,  $= 0$  when it is a boundary point.

The same type of reasoning also yields a more general relation,

$$q_n(x, y) = \sum_{z \in V \setminus V_0} q_{n_1}(x, z) q_{n-n_1}(z, y) \text{ whenever } n_1 < n \quad (4.2.97)$$

where the sum now extends over all interior vertices. This relation follows from the simple observation that whatever reaches  $y$  from  $x$  in  $n$  steps has to be at some interior vertex at the time  $n_1 < n$  whence it arrives at  $y$  after  $n - n_1$  further steps. Obviously, (4.2.96) is a special case of (4.2.97), corresponding to  $n_1 = n - 1$ .

Returning to (4.2.96), we see that the partial sums

$$v_n(x, y) := \sum_{\nu=0}^n q_\nu(x, y) \quad (4.2.98)$$

satisfy

$$v_n(x, y) = \begin{cases} \sum_{y' \sim y} \frac{1}{n_{y'}} v_{n-1}(x, y') & \text{for } x \neq y \\ 1 + \sum_{y' \sim y} \frac{1}{n_{y'}} v_{n-1}(x, y') & \text{for } x = y \end{cases} \quad (4.2.99)$$

We define the operator  $\Delta'$  by

$$\Delta' f(y) := \sum_{y' \sim y} \frac{1}{n_{y'}} f(y') - f(y). \quad (4.2.100)$$

When our graph is regular in the sense that all vertices  $y$  have the same degree  $n_y = k$ , say, then  $\Delta' = \Delta$ . For other graphs, the two operators are obviously different.

From (4.2.96), we infer

$$q_{n+1}(x, y) - q_n(x, y) = \Delta'_y q_n(x, y) \quad (4.2.101)$$

and  $q_n(x, y) = 0$  when  $y$  is a boundary point. This is a discrete heat equation, with time step 1 and  $\Delta'$  in place of  $\Delta$ .

Similarly, from (4.2.99), (4.2.98), we infer

$$\Delta'_y v_n(x, y) = \begin{cases} q_n(x, y) & \text{for } x \neq y \\ q_n(x, y) - 1 & \text{for } x = y \end{cases} \quad (4.2.102)$$

and  $v_n(x, y) = 0$  when  $y$  is a boundary point. Again, we can rewrite this as a heat type equation

$$\Delta'_y v_n(x, y) = \begin{cases} v_{n+1}(x, y) - v_n(x, y) & \text{for } x \neq y \\ v_{n+1}(x, y) - v_n(x, y) - 1 & \text{for } x = y \end{cases} \quad (4.2.103)$$

Since we already know that  $q_n(x, y) \rightarrow 0$  for  $n \rightarrow \infty$ ,  $v_n(x, y)$  converges to the solution  $v$  of

$$\Delta'_y v(x, y) = \begin{cases} 0 & \text{for } x \neq y \\ -1 & \text{for } x = y \end{cases} \quad (4.2.104)$$

with boundary values 0.

When we vary  $x$  in place of  $y$ , we come up with the Laplacian  $\Delta$  in place of  $\Delta'$ . Indeed, we have

$$q_n(x, y) = \frac{1}{n_x} \sum_{x' \sim x} q_{n-1}(x', y) \quad (4.2.105)$$

because any random path that goes from  $x$  to  $y$  in  $n$  steps has to pass through one of the neighbors of  $x$  in the first step with equal probability  $\frac{1}{n_x}$ . (This is again a special case of (4.2.97), this time for  $n_1 = 1$ .) Therefore, we obtain the discrete heat equation

$$q_{n+1}(x, y) - q_n(x, y) = \Delta_x q_n(x, y). \quad (4.2.106)$$

As before, from (4.2.105), we conclude

$$v_n(x, y) = \begin{cases} \frac{1}{n_x} \sum_{x' \sim x} v_{n-1}(x', y) & \text{for } x \neq y \\ 1 + \frac{1}{n_x} \sum_{x' \sim x} v_{n-1}(x', y) & \text{for } x = y \end{cases} \quad (4.2.107)$$

From (4.2.99), (4.2.107), we infer

$$\Delta_x v_n(x, y) = \begin{cases} q_n(x, y) & \text{for } x \neq y \\ q_n(x, y) - 1 & \text{for } x = y \end{cases} \quad (4.2.108)$$

and  $v_n(x, y) = 0$  when  $x$  is a boundary point or, equivalently,

$$\Delta_x v_n(x, y) = \begin{cases} v_{n+1}(x, y) - v_n(x, y) & \text{for } x \neq y \\ v_{n+1}(x, y) - v_n(x, y) - 1 & \text{for } x = y. \end{cases} \quad (4.2.109)$$

Since we already know that  $q_n(x, y) \rightarrow 0$  for  $n \rightarrow \infty$ ,  $v_n(x, y)$  converges to the solution  $v$  of

$$\Delta_x v(x, y) = \begin{cases} 0 & \text{for } x \neq y \\ -1 & \text{for } x = y \end{cases} \quad (4.2.110)$$

with boundary values 0. Up to the normalization factor and the minus sign in (4.2.104), this solution is the Green function as defined in (4.1.32). In particular, the solution of the Poisson problem

$$\Delta u(x) = g(x) \text{ for } x \in S \setminus V_0 \quad (4.2.111)$$

$$u(x) = 0 \text{ for } x \in V_0 \quad (4.2.112)$$

is given by

$$u(x) = - \sum_y v(x, y) g(y), \quad (4.2.113)$$

that is, the negative of the expected sum of  $g$  along the random walk starting at  $x$  until it reaches the boundary.

Finally,  $w$  from (4.2.93) satisfies

$$w(x, z) = \frac{1}{n_x} \sum_{x' \sim x} w(x', z) \quad (4.2.114)$$

because any path from  $x$  to  $z$  has to pass through one of the neighbors of  $x$ . This means

$$\Delta_x w(x, z) = 0. \quad (4.2.115)$$

For two different boundary points  $z_1, z_2$ , we have  $w(z_1, z_2) = 0$ , and  $w(z, z) = 1$ . Thus,  $w(x, z)$  as a function of  $x$  solves the Dirichlet problem (see (4.1.29)) with those boundary values. In words: the probability as a function of the starting point  $x$  of the random walk for being absorbed at the boundary point  $z$  is the

harmonic function  $u$  on the graph with boundary values  $u(y, z) = \delta(y, z)$ . For general boundary values  $f(z)$  for  $z \in V_0$ , the solution of the Dirichlet problem is

$$u(x) = \sum_z w(x, z)f(z). \quad (4.2.116)$$

According to our above interpretation of  $w(x, \cdot)$  as the probability distribution for the exit point of the random walk starting at  $x$ , we can express (4.2.116) as follows: The solution  $u(x)$  at the point  $x$  of the Dirichlet problem with boundary values  $f$  is the expected value of  $f$  at the exit point for the random walk starting at  $x$ ,

$$u(x) = E(f(w(x, \cdot))). \quad (4.2.117)$$

Obviously, the position  $X_n$  of the random walker on our graph constitutes a random process in the sense of definition 3.2.1. It also satisfies the Markov property of definition 3.2.2 because the probability distribution for the position  $X_{n+1}$  depends only on the location  $X_n = x$  at time  $n$ , but is independent of earlier positions when given that position at time  $n$ .

We now briefly consider the case without boundary, that is,  $V_0 = \emptyset$ . The transition probabilities for  $X_n$  are independent of  $n$  and given

$$P(x, y) := p(X_{n+1} = y | x_n = x) = \frac{1}{n_x}. \quad (4.2.118)$$

In the above, we have considered the initial distribution  $f_0(y) = \delta(x, y)$  (the random walker always started at the point  $x$ ), but we can obviously consider any initial distribution  $f_0$  with  $\sum_y f_0(y) = 1$ . Given an initial distribution  $f_0$ , the distribution  $f_n$  at time  $n$  then is  $f_0 P^n$  where  $f_0$  is considered as a row vector, that is,

$$f_n(x_n) = f_0(x_0)P(x_0, x_1) \dots P(x_{n-1}, x_n). \quad (4.2.119)$$

A distribution  $\pi$  is called stationary if

$$\pi P = \pi, \text{ that is, } \pi(y) = \pi(x)P(x, y). \quad (4.2.120)$$

The random walk is called ergodic if there exists a unique stationary distribution  $\pi$  with

$$\lim_{n \rightarrow \infty} f_0 P^n = \pi \quad (4.2.121)$$

for every initial distribution  $f_0$ . The process is ergodic iff it is irreducible, i.e., for every  $x, y \in V$  there exists some  $n$  with  $P^n(x, y) > 0$ , and aperiodic, i.e., the greatest common divisor of the  $n$  with  $P^n(x, y) > 0$  is 1. The first condition is equivalent to the graph  $\Gamma$  being connected, or in terms of eigenvalues  $\lambda_1 > 0$ , while the second one is equivalent to  $\Gamma$  being not bipartite, that is, the largest eigenvalue  $\lambda_K < 2$ , see (2.1.37) and (2.1.39), resp., in 2.1.3.

We now want to turn to the continuous case (a good reference is [15]). Our heuristic strategy consists in taking a regular lattice as our graph and pass to the continuum limit. This means that we consider the lattice of points

$\{h(n_1, \dots, n_d) : n_1, \dots, n_d \in \mathbb{Z}\}$  for  $h > 0$  which we want to let tend to 0. Thus, our random walker on this lattice when at the lattice point  $z_n$  at time  $n$  moves to one of its  $2d$  lattice neighbors with equal probability  $\frac{1}{2d}$ . Then the random variable  $Z_n = (Z_n^1, \dots, Z_n^d)$  describing the position of the random walker at time  $n$  satisfies

$$Z_n^j - Z_0^j = \sum_{i=1}^n X_i \quad (4.2.122)$$

where the  $X_i$  are independent identically distributed random variables with probabilities

$$p(X_i = h) = \frac{1}{2d}, \quad p(X_i = -h) = \frac{1}{2d}, \quad p(X_i = 0) = \frac{d-1}{d} \quad (4.2.123)$$

where  $X_i = 0$  corresponds to the case where the random step is taking in a direction that is not the  $i$ th coordinate direction. The  $X_i$  all have expectation value 0 and standard deviation  $\frac{1}{d}h^2$ . The  $Z_n^j - Z_0^j$  then also have expectation value 0, and their standard deviation is  $\frac{n}{d}h^2$ , by Lemma 3.1.2. By the central limit theorem 3.1.2, for  $n \rightarrow \infty$ ,  $Z_n^j - Z_0^j$  approaches the Gaussian distribution  $N(0, \frac{1}{d}nh^2) = \frac{1}{\sqrt{2\pi\frac{1}{d}nh^2}} \exp(-\frac{x^2}{2\frac{1}{d}nh^2})$ . We also want to let the size of the time step tend to 0, to compensate for the factor  $h^2$  going to 0. That is, we let the random walker move at times  $\tau, 2\tau, 3\tau, \dots$ . Then at time  $t = m\tau$ , it has jumped  $m$  times, and the corresponding position  $Z^j(t) - Z^j(0)$  is distributed according to  $N(0, \frac{1}{d}\frac{t}{\tau}h^2)$ . In order to have a nontrivial limit for positive finite  $t$ , we then let  $h$  and  $\tau$  tend to 0 in such a manner that  $\frac{h^2}{d\tau} =: \mu^2$  is a positive constant. Thus, in the limit,  $Z^j(t) - Z^j(0)$  is distributed according to  $N(0, t\mu^2)$ . The limiting process  $X(t) = (X^1(t), \dots, X^d(t))$  – whose existence one needs to prove – is called the Wiener process or Brownian motion. The components  $X^j(t)$  are independent and identically distributed (this is a consequence of the homogeneity of the lattice and the fact that the random walker was moving in each direction with the same probability).  $X^j(t+s) - X^j(t)$  is distributed according to  $N(0, \mu^2 s)$ . In particular, this does not depend on  $t$ . Moreover,  $X(t_1) - X(s_1)$  and  $X(t_2) - X(s_2)$  are independent whenever  $s_1 < t_1 < s_2 < t_2$  – one says that  $X$  has independent increments (cf. the corresponding notion introduced above for point processes). Finally, the typical path  $X(t), t \geq 0$  is continuous (but nowhere differentiable).

Again, when we have a bounded domain  $\Omega \subset \mathbb{R}^d$  and prescribe continuous boundary values  $f$  on  $\partial\Omega$ , the Dirichlet problem (cf. (4.1.29))

$$\Delta u(x) = 0 \quad \text{for } x \in \Omega \quad (4.2.124)$$

$$u(z) = f(z) \quad \text{for } z \in \partial\Omega \quad (4.2.125)$$

can be solved by Brownian motion: the (unique) solution  $u(x)$  at the point  $x$  of the Dirichlet problem with boundary values  $f$  is the expected value of  $f$  at the exit point for the random walk starting at  $x$ ,

$$u(x) = E(f(W(x, \cdot))) \quad (4.2.126)$$

where the random variable  $W(x, \cdot)$  encodes the exit point from  $\Omega$  for the random walk starting at  $x$ . There is one technical issue here, namely about attaining the boundary values, that is for which points  $z \in \partial\Omega$  we have

$$\lim_{x \rightarrow z, x \in \Omega} E(f(W(x, z))) = f(z). \quad (4.2.127)$$

The points in  $\partial\Omega$  satisfying this condition are called regular. They can be characterized in potential theoretic terms. In particular, this does not depend on the (continuous) function  $f$ , but only on the geometry of the domain  $\Omega$ . Not every point is regular, however. For example, for  $d \geq 2$ , isolated boundary points are not regular (because they constitute removable singularities for harmonic functions). Here, we do not intend to go into this issue in more detail.

Likewise, up to the minus sign, the Green function is given by the solution  $v(x, y)$  of the analogue of (4.2.104) or (4.2.110). In particular, the Poisson problem

$$\Delta u(x) = g(x) \text{ for } x \in \Omega \quad (4.2.128)$$

$$u(z) = 0 \text{ for } z \in \partial\Omega \quad (4.2.129)$$

is given by

$$u(x) = - \int v(x, y)g(y)dy. \quad (4.2.130)$$

Here,  $v(x, y)$  is the negative of the Green function, that is, the solution of

$$\Delta_x v(x, y) = -\delta(x, y) \text{ for } x \in \Omega \quad (4.2.131)$$

$$v(x, y) = 0 \text{ for } x \in \partial\Omega, \quad (4.2.132)$$

in analogy to (4.2.110), (4.2.111), (4.2.113).

Similarly, the following interpretation is carried over from the discrete case: For  $A \subset \Omega$ ,

$$v(x, A) := \int_A v(x, y)dy \quad (4.2.133)$$

is the expected amount of time the random walk starting at  $x$  spends in  $A$  before exiting from  $\Omega$ . In probabilistic terminology, (4.2.130) is also expressed as

$$u(x) = -E\left(\int_0^{\tau_\Omega} g(X_x(t))dt\right) \quad (4.2.134)$$

where  $X_x(t)$  is Brownian motion starting at  $x$  and  $\tau_\Omega$  is its expected exit time from  $\Omega$ . In words: the solution  $u(x)$  at  $x$  of the Poisson problem for  $g$  is given by the negative of the expected integral of  $g$  over a random path starting at  $x$  until it exits from  $\Omega$ . In particular, we may put  $g = -1$ . Then (4.2.134) becomes

$$u(x) = E(\tau_\Omega), \quad (4.2.135)$$

the expected exit time of the random walk starting at  $x$ . Thus, this expected exit time is the solution of

$$\Delta u(x) = -1 \text{ in } \Omega, \quad u(y) = 0 \text{ for } y \in \partial\Omega. \quad (4.2.136)$$

We return to the probability density

$$P(y, t|x, s) := p(X(t) = y|X(s) = x) = \frac{1}{\sqrt{2\pi(t-s)}} \exp\left(-\frac{(y-x)^2}{2(t-s)}\right) \quad (4.2.137)$$

for  $t > s$ . This probability density satisfies

$$\frac{\partial P}{\partial t} = \frac{1}{2} \Delta_y P \quad (4.2.138)$$

and

$$\frac{\partial P}{\partial s} = -\frac{1}{2} \Delta_x P \quad (4.2.139)$$

(4.2.138) is called the forward diffusion or Kolmogorov equation, (4.2.139) the backward one. (4.2.138) is also called the Fokker-Planck equation. The interpretation is that the probability density of a stochastic process (here the Wiener process or Brownian motion) satisfies a *deterministic* differential equation.

(4.2.138) is the continuous analogue of (4.2.101). We obtain the Laplace operator here because the lattice that we used for our approximation scheme was regular as all vertices had the same degree  $2d$ .

We also have an analogue of (4.2.97),

$$P(y, t|x, s) = \int_z P(y, t|z, s+\tau)P(z, s+\tau|x, s) dz \text{ for } 0 < \tau < t-s. \quad (4.2.140)$$

Again, the reason for this relation is that whatever arrives at time  $t$  at  $y$ , originating from  $x$  at time  $s$  has to be at some point  $z$  at the intermediate time  $s+\tau$  whence it reaches  $y$  after the further time  $t-\tau$ . (4.2.140) is called the Chapman-Kolmogorov equation. Of course, (4.2.140) can also be derived by a direct computation with Gaussian kernels, on the basis of (4.2.137), but our more abstract derivation is simpler and more insightful. In any case, we again see the ubiquity of Gaussian kernels. By the central limit theorem, our rescaling scheme for the random walk on a lattice produced a Gaussian transition kernel which in turn governs the standard heat equation.

This can also be coupled with a deterministic drift. We consider a general dynamical rule of the form

$$\frac{dy}{dt} = F(y(t), t), \text{ for } y \in \mathbb{R}^d \quad (4.2.141)$$

The continuity equation for the density of  $y$  then is

$$\frac{\partial}{\partial t} p(y, t) = -\sum_{i=1}^d \frac{\partial}{\partial y^i} (F^i(y, t)p(y, t)). \quad (4.2.142)$$

We now take the sum of Brownian motion and a deterministic dynamics of the form (4.2.141). We write this formally as

$$\frac{dy}{dt} = F(y(t), t) + \eta, \quad (4.2.143)$$



where  $\eta$  is the formal derivative of Brownian motion (which is represented by white noise, but we do not explain this here; see e.g. [27, 30, 22]). (The equation (4.2.143) is called the Langevin equation.) By linear superposition of (4.2.138) and (4.2.142), the corresponding density satisfies the Fokker-Planck equation

$$\frac{\partial}{\partial t}p(y, t) = \frac{1}{2}\Delta p(y, t) - \sum_{i=1}^d \frac{\partial}{\partial y^i}(F^i(y, t)p(y, t)). \quad (4.2.144)$$

This issue will be taken up again in 4.5 below.

## 4.3 Dynamical systems

### 4.3.1 Systems of ordinary differential equations

A general reference for this section is [22]. Let  $f = (f^1, \dots, f^n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be of class  $C^1$ . We consider the system of first order ordinary differential equations (ODEs)<sup>2</sup>

$$\dot{x}^i(t) = f^i(x^1(t), \dots, x^n(t)) \text{ for } i = 1, \dots, n, \quad (4.3.145)$$

with  $\dot{x}^i = \frac{d}{dt}x^i$ .  $t$  is considered to be the time, and  $x(t) = (x^1(t), \dots, x^n(t))$  then is the state of the system at time  $t$ . One usually prescribes initial values  $x_0 = x(0)$  and looks for a solution  $x(t), t \in \mathbb{R}(t \geq 0)$ .  $\{x(t) : t \geq 0\}$  is called the **orbit** of  $x_0$ .

(4.3.145) is a so-called **autonomous** system because  $f$  does not depend explicitly on  $t$  (but implicitly through the dependence of  $x$  on  $t$ ).<sup>3</sup> The important point about autonomous systems is that they are invariant under time shifts. This means that, if we consider the solution of (4.3.145)<sup>4</sup>  $x_1(t)$  with initial values  $x_1(t_1) = \xi$  and the solution  $x_2(t)$  with the same initial values, but starting at time  $t_2$ , that is,  $x_2(t_2) = \xi$ , then for all  $t \geq t_2$ ,  $x_2(t) = x_1(t + t_1 - t_2)$ . In other

<sup>2</sup>Higher order systems of ODEs can be reduced to systems of first order by introducing additional auxiliary variables.

<sup>3</sup>One may also consider non-autonomous systems,

$$\dot{x}^i(t) = \phi^i(t, x^1(t), \dots, x^n(t)) \text{ for } i = 1, \dots, n,$$

with an explicit dependence on  $t$ , but such systems can be converted into an autonomous form by introducing a new dependent variable  $x^{n+1}$  to obtain the equation

$$\dot{x}^{n+1}(t) = f^{n+1}(x^1(t), \dots, x^n(t), x^{n+1}(t)) \equiv 1.$$

This may turn linear (non-autonomous) equations into non-linear (autonomous) ones; e.g.

$$\dot{x} = \cos(\beta t)$$

becomes

$$\begin{aligned} \dot{x}^1 &= \cos(\beta x^2) \\ \dot{x}^2 &= 1 \end{aligned}$$

where the dependent variable  $x^2$  enters the r.h.s. non-linearly.

<sup>4</sup>assuming that there exists a unique solution, see below

words, the behavior of the solution (obviously) depends on the initial values, that is, where or how it starts, but not on the starting time, that is, when it starts.

The theorem of Picard-Lindelöf yields the short-time existence of solutions:

**Theorem 4.3.1.** *Suppose that  $f$  is Lipschitz continuous, that is, there exists some constant  $L$  with*

$$|f(x_1) - f(x_2)| \leq L|x_1 - x_2| \quad (4.3.146)$$

for all  $x_1, x_2 \in \mathbb{R}^n$ .

For every initial state  $x_0$ , the solution  $x(t)$  of the system(4.3.145) then exists on some time interval, that is, for

$$-T < t < T, \text{ for some } T > 0.$$

*This solution is unique.*

The solution need not exist for all times, that is, the maximal such  $T$  may be finite. That maximal  $T$  in general depends on the initial values  $x_0$ . We shall see examples shortly.

An easy, but important consequence of the uniqueness part of the Picard-Lindelöf theorem is that orbits of an autonomous system (4.3.145) cannot intersect or merge. Namely, when at some point  $x_0$  two orbits came together, then there would exist two different solutions (in forward or backward time) with initial values given by  $x_0$ .

Since we are imposing no restrictions on  $f$  apart from a rather mild smoothness assumption, the behavior of the solutions of systems of ODEs can be rather diverse, and one cannot expect a useful classification. It is more insightful to study certain dynamical motives, that is, qualitative types of behavior of solutions. We start with the case  $n = 1$  and write  $f$  in place of  $f^1$ , and likewise for  $x$ . Thus, we look at the scalar equation

$$\dot{x} = f(x) \text{ for } t \geq 0, \text{ with initial condition } x(0) = x_0. \quad (4.3.147)$$

Clearly, there are the simple linear equations, like

$$\dot{x} = 0 \quad (4.3.148)$$

whose solution is constant,  $x(t) = x_0$ , or

$$\dot{x} = b \quad (4.3.149)$$

whose solution is linear,  $x(t) = x_0 + bt$ , or

$$\dot{x} = cx \quad (4.3.150)$$

whose solution is  $x(t) = x_0 e^{ct}$ . The latter equation, for  $c > 0$ , already exhibits the important phenomenon that solutions of ODEs can amplify differences over

time, that is, when we have two solutions  $x_1, x_2$  with different initial values  $x_i(0) = x_{0,i}$ , then  $|x_1(t) - x_2(t)| = |x_{0,1} - x_{0,2}|e^{ct}$  grows exponentially. Of course, exponential growth cannot be sustained for a long time. Thus, in many models, one introduces a carrying capacity and considers

$$\dot{x} = cx(m - x) \quad (4.3.151)$$

for  $c, m > 0$ . This equation is called the logistic, Verhulst, or Fisher equation. Below, we shall often consider this equation as an example, usually for  $c = m =$ , that is,

$$\dot{x}x(1 - x). \quad (4.3.152)$$

In (4.3.151), for initial values  $0 \leq x(0) \leq m$ , the solution is bounded and stays in that same interval,  $0 \leq x(t) \leq m$  for all  $t \geq 0$ . In that case,  $x(t)$  is monotonically increasing, with  $\lim_{t \rightarrow \infty} x(t) = m$ . When  $x(0) > m$ , the solution monotonically decays towards the asymptotic value  $m$ .  $x = m$  and  $x = 0$  are both fixed points, that is, when  $x_0 = m$  or  $0$ , then  $\dot{x}(t) = 0$  for all times, and the solution will stay constant. When the initial values are negative, however, the solution diverges to  $-\infty$  in finite time. In particular, we here see the phenomenon that a solution need not exist for all positive times; the simplest example of this is perhaps

$$\dot{x} = x^2, \quad (4.3.153)$$

with the solution

$$x = \left(\frac{1}{x(0)} - t\right)^{-1} \quad (4.3.154)$$

which when  $x(0)$  is positive becomes infinite in finite time. In fact, the blow-up occurs at  $t = \frac{1}{x(0)}$ . In contrast, when  $x(0) < 0$ , the solution exists for all time, with  $\lim_{t \rightarrow \infty} x(t) = 0$ . When  $x(0) = 0$ , then  $x(t) = 0$  for all  $t$ . Thus, the fixed point at  $x = 0$  separates two different qualitative regimes for the solution of our differential equation.

In general, if  $f(x_*) = 0$  for  $i = 1, \dots, d$  then  $x_*$  is a fixed point for our equation (4.3.147), that is, when  $x_0 = x_*$ , then  $x(t) = x_*$  for all  $t$ .

Our differential equation (4.3.147) may have several fixed points  $x_1 < x_2 < \dots < x_m$ . (We assume here for simplicity that there are only finitely many fixed points. The case of infinitely many fixed points does not lead to substantially new phenomena as the reader will easily check.) If there is no further fixed point between  $x_k$  and  $x_{k+1}$ , then  $f(x)$  cannot have a zero for  $x_k < x < x_{k+1}$  and therefore must have a definite sign there. When this sign is positive, then for  $x_k < x(0) < x_{k+1}$ , the solution  $x(t)$  of (4.3.147) – which exists for all  $t$  – satisfies  $\lim_{t \rightarrow \infty} x(t) = x_{k+1}$ . Similarly, when  $f(x) < 0$  for  $x_k < x < x_{k+1}$ , the solution with initial values in that interval satisfies  $\lim_{t \rightarrow \infty} x(t) = x_k$ . In particular, the fixed point  $x_k$  is attracting when  $f(x) > 0$  for  $x_{k-1} < x < x_k$  and  $f(x) < 0$  for  $x_k < x < x_{k+1}$ . It is repelling when both signs are reversed. The fixed point  $0$  for (4.3.153) is neither attracting nor repelling, because  $f(x)$  does not change its sign there. When  $x_m$  is the largest fixed point, then either  $f(x) > 0$  for  $x > x_m$  in which case the solution could possibly blow up in finite

time, or  $f(x) < 0$  for  $x > x_m$  in which case the solution monotonically decays to  $x_m$  for initial values  $x(0) > x_m$ . The analogous situation holds when the initial values are smaller than the smallest fixed point.

We can also formulate the following easy global existence result:

**Theorem 4.3.2.** *We consider (4.3.147),*

$$\dot{x} = f(x) \text{ for } t \geq 0, \text{ with initial condition } x(0) = x_0 \quad (4.3.155)$$

and assume that there exist numbers  $m < M$  with

$$f(m) > 0 \text{ and } f(M) < 0 \quad (4.3.156)$$

and

$$m \leq x_0 \leq M. \quad (4.3.157)$$

Then (for a Lipschitz continuous  $f$  as in (4.3.146)), the solution of (4.3.155) exists for all  $t \geq 0$ .

*Proof.* The key observation is that the solution  $x(t)$  has to stay bounded as long as it exists. Whenever it comes near the upper bound  $M$ , then  $\dot{x} = f(x)$  becomes negative by (4.3.181), and therefore  $x(t)$ , and when it comes near the lower bound  $m$ , it increases for the same reason. Therefore, we shall have  $m \leq x(t) \leq M$  for all  $t$  for which the solution exists. By the theorem of Picard-Lindelöf, we can then find some  $T > 0$  such that for each  $t_0$  up to which the solution exists, the solution of  $\dot{y}(t) = f(y(t))$  with  $y(0) = x(t_0)$  exists for  $0 \leq tT$ . The important point here is that  $T$  here does not depend on  $x(t_0)$  because the latter is confined in the compact interval  $[m, M]$ .  $x(t) = y(t - t_0)$  then is the solution of (4.3.155) on the interval  $[t_0, t_0 + T]$ . This implies that the solution has to exist for all time (negative times, although not really our concern, are handled by the same argument).  $\square$

When the r.h.s. of (4.3.147) depends on some parameter  $\lambda$ , that is, we look at

$$\dot{x} = f(x, \lambda) \quad (4.3.158)$$

then we expect a bifurcation, that is, a qualitative change of behavior of the solutions at those parameter values  $\lambda = \lambda_0$  where the number of solutions  $x_k$  of

$$f(x, \lambda) = 0 \quad (4.3.159)$$

changes. For example, for

$$\dot{x} = x^2 + \lambda \quad (4.3.160)$$

$\lambda = 0$  is such a bifurcation point. For  $\lambda > 0$ , there is no fixed points, for  $\lambda = 0$  there is one, namely 0, and for  $\lambda < 0$ , there are two,  $x = \pm\sqrt{-\lambda}$ . Here, we already see the important principle that at generic (that is, typical) bifurcations, fixed always arise or disappear in pairs.

For

$$\dot{x} = \lambda x - x^3 \quad (4.3.161)$$

$\lambda = 0$  is again a bifurcation point. For  $\lambda \leq 0$ ,  $x = 0$  is the only fixed point whereas for  $\lambda > 0$ , we have additional ones at  $x = \pm\sqrt{\lambda}$ . The latter ones are attracting whereas 0 is repelling for  $\lambda > 0$ , but attracting for  $\lambda \leq 0$ .

The preceding already summarizes the main qualitative results about single ODEs of first order. For  $n > 1$ , the behavior of solutions of (4.3.145) can become richer and more interesting. When we move to dimension  $n = 2$ , two new phenomena emerge:

- saddle type fixed points in addition to attracting and repelling ones
- closed periodic orbits.

This is best understood within the context of some wider principles for the analysis of dynamical systems:

1. identify the compact orbits (and perhaps other invariant sets) of the dynamics,
2. linearize about them and
3. investigate their stability.

We now elaborate these points. The simplest case of an invariant set is a fixed point, that is, a point  $x_*$  for which the rhs of (4.3.145) vanishes, that is,

$$f^i(x_*^1, \dots, x_*^n) = 0. \quad (4.3.162)$$

Typically, the investigation of a dynamical system starts with the identification of these fixed points. We may assume  $x_* = 0$  and study the linearized system

$$\dot{x}(t) = Lx, \text{ with } L = \left( \frac{\partial f^i}{\partial x^j}(x_*) \right)_{i,j=1,\dots,n}. \quad (4.3.163)$$

We consider the case  $n = 2$  because the above principles already become clear there. The matrix  $L$  has either two real eigenvalues or two conjugate complex ones. When it can be diagonalized with two real eigenvalues  $\lambda_1$  and  $\lambda_2$ , the after a linear change of coordinates, our linearized system becomes

$$\begin{aligned} \dot{x}^1(t) &= \lambda_1 x^1(t) \\ \dot{x}^2(t) &= \lambda_2 x^2(t), \end{aligned} \quad (4.3.164)$$

the solution of which obviously is

$$\begin{aligned} x^1(t) &= e^{\lambda_1 t} x^1(0) \\ x^2(t) &= e^{\lambda_2 t} x^2(0). \end{aligned} \quad (4.3.165)$$

If both eigenvalues are negative, then  $x(t)$  converges to the fixed point  $x_*(=0)$  exponentially while, in the case where both are positive,  $x(t)$  exponentially expands.

In the first, attracting, case,  $x_* = 0$  is called a **node** or **sink**, and it is a stable fixed point for  $t \rightarrow \infty$ , whereas in the second, repelling, case, called a **source**, it is unstable for  $t \rightarrow \infty$ . If the two eigenvalues have different signs, say  $\lambda_2 < 0 < \lambda_1$ , then the fixed point  $x_* = 0$  is neither stable nor unstable. In fact, any initial point on the  $x^2$ -axis converges to 0, while all other initial points diverge under the flow. This is called a **saddle**. When one of the eigenvalues vanishes, the picture can get more complicated, and the behaviour of the linearized system may be different from the original one. Actually, this is already seen in the one-dimensional example

$$\dot{x} = x^2 \quad (4.3.166)$$

the linearization of which at 0 is

$$\dot{x} = 0. \quad (4.3.167)$$

When, in contrast to the preceding cases,  $L$  has two complex conjugate eigenvalues  $\lambda \pm i\mu$ , then, after a linear change of coordinates again, we get the system

$$\dot{x}(t) = \begin{pmatrix} \lambda & \mu \\ -\mu & \lambda \end{pmatrix} x(t), \quad (4.3.168)$$

the solution of which is

$$x(t) = e^{\lambda t} \begin{pmatrix} \cos \mu t & \sin \mu t \\ -\sin \mu t & \cos \mu t \end{pmatrix} x(0). \quad (4.3.169)$$

For  $\lambda \neq 0$ ,  $x(t)$  moves on a spiral (clockwise or counterclockwise, depending on the sign of  $\mu$ ), exponentially towards 0 for  $\lambda < 0$ , exponentially expanding for  $\lambda > 0$ . The case  $\lambda = 0$  seems intermediate as the solution then moves on a circle about 0. This case is important for two reasons: it is the first non-trivial example of a compact orbit, and thus an invariant set, other than a point. Secondly, this case  $\lambda = 0$  is different from all previous cases, insofar as it is not structurally stable. This means that an arbitrarily small variation of  $\lambda = 0$  changes the qualitative behavior of the system. Even worse, even the qualitative behavior of the linearized system near the fixed point 0 is no longer the same as that of the original system in case  $\lambda = 0$ . In a certain sense, these two phenomena are related as we shall try to explain soon. Before doing that, we formulate a general

**Definition 4.3.1.** A fixed point  $x_*$  of (4.3.145) is called hyperbolic if all eigenvalues of the linearized system have nonvanishing real part.

The qualitative dynamical behaviour near a hyperbolic fixed point is structurally stable in the sense that it is not affected by sufficiently small perturbations or parameter variations (like the eigenvalues of the linearized system), and that it is the same as the one of the linearized system – in fact, the difference between the original and the linearized system is such a small perturbation that

does not change the qualitative behaviour.

We now look into a non-hyperbolic situation and consider the linear system

$$\dot{x} = y + \lambda x \quad (4.3.170)$$

$$\dot{y} = -x + \lambda y. \quad (4.3.171)$$

The eigenvalues are  $\lambda \pm i$ , with imaginary part  $\neq 0$ , and real part  $= 0$  for  $\lambda = 0$ . We may consider  $\lambda$  as a bifurcation parameter, and  $\lambda = 0$  as a bifurcation value where the qualitative behaviour of the solution changes. Here, a pair of complex conjugate nonzero eigenvalues crosses the imaginary axis. This is the characteristic criterion for the so-called Hopf bifurcation as we shall now explain.

In the linear system, at  $\lambda = 0$  all orbits are periodic, namely circles, about  $(0, 0)$ , while for  $\lambda \neq 0$  there is no periodic orbit at all.

(4.3.170) is the linearization at  $(0, 0)$  of

$$\dot{x} = y - x(x^2 + y^2 - \lambda) \quad (4.3.172)$$

$$\dot{y} = -x - y(x^2 + y^2 - \lambda) \quad (4.3.173)$$

For this system,  $(0, 0)$  is a fixed point for all parameter values  $\lambda$ . For  $\lambda \neq 0$ , the situation is hyperbolic and therefore qualitatively the same as in the linearized system. For  $\lambda < 0$ , the fixed point is globally exponentially attracting. While this can be deduced from general principles, we can, of course, also see it directly, and this gives us the opportunity to introduce another useful tool, a Lyapunov function which by definition is a function that is strictly decreasing along every flow line. Here, such a Lyapunov function is given by  $\log(x^2 + y^2)$  since we have

$$\frac{d}{dt} \log(x^2 + y^2) = 2(-x^2 - y^2 + \lambda) \leq 2\lambda < 0. \quad (4.3.174)$$

Thus,  $\log(x^2 + y^2)$  decreases along every flow line, and then so does  $x^2 + y^2$ , and therefore each flow line has to lead to  $(0, 0)$ . As already pointed out, this is a structurally stable situation that is invariant under small perturbations of  $\lambda$ .

For  $\lambda = 0$ ,  $(0, 0)$  is still globally attracting, but no longer exponentially so. We still have

$$\frac{d}{dt} \log(x^2 + y^2) < 0 \text{ for } (x, y) \neq (0, 0), \quad (4.3.175)$$

but this expression is no longer bounded away from 0. Thus, we see again that the situation at  $\lambda = 0$  is not structurally stable.

For  $\lambda > 0$ , while the situation near  $(0, 0)$  is again structurally stable, an interesting global phenomenon emerges away from  $(0, 0)$ .  $(0, 0)$  is repelling, and there exists a periodic orbit  $x^2 + y^2 = \lambda$  that is attracting. To understand this, we consider our Lyapunov function:

$$\frac{d}{dt} \log(x^2 + y^2) \begin{cases} > 0 & \text{for } x^2 + y^2 < \lambda \\ = 0 & \text{for } x^2 + y^2 = \lambda \\ < 0 & \text{for } x^2 + y^2 > \lambda. \end{cases} \quad (4.3.176)$$

Thus, when we are on the circle  $x^2 + y^2 = \lambda$ , we stay there and since,  $\dot{x}$  and  $\dot{y}$  do not vanish there, it is a nontrivial periodic orbit. When we are outside or inside that circle, we move towards it.

We thus obtain a family, depending on  $\lambda$ , of periodic orbits that emerge from the fixed point at the transition from  $\lambda = 0$  to  $\lambda > 0$ . This family of periodic orbits represents a structurally stable bifurcation, that is, such a family remains under perturbations of the above system.

In contrast to this behaviour, in the linear system, at  $\lambda = 0$  all orbits are periodic, namely circles, about  $(0, 0)$ , while for  $\lambda \neq 0$  there is no periodic orbit at all. Thus, here the whole family of periodic orbits is concentrated at a single parameter value, while when the linear system is perturbed by a higher order term, that family gets distributed among different parameter values. The situation at  $\lambda = 0$  itself is not structurally stable while the behaviour of the whole family is, namely the emergence of a family of periodic orbits at the transition from an attracting to a repelling fixed point.

We next consider another system with the same linearization (4.3.170), (4.3.171) as the preceding one, (4.3.172), (4.3.173),

$$\dot{x} = y - x((x^2 + y^2)^2 - 2(x^2 + y^2) - \lambda) \quad (4.3.177)$$

$$\dot{y} = -x - y((x^2 + y^2)^2 - 2(x^2 + y^2) - \lambda) \quad (4.3.178)$$

depending on a real parameter  $\lambda$  as before. We now have

$$\frac{d}{dt} \log(x^2 + y^2) = 2(-(x^2 + y^2)^2 + 2(x^2 + y^2) + \lambda). \quad (4.3.179)$$

This becomes 0 when

$$x^2 + y^2 = 1 \pm \sqrt{1 + \lambda}.$$

Thus, whenever this value is real and nonnegative, we obtain that  $x^2 + y^2$  remains constant along a solution, that is, the orbit is a circle. When  $\lambda$  is smaller than  $-1$ , no such solution exists. For  $\lambda = -1$ , we find precisely one solution whereas, for  $-1 < \lambda < 0$ , we obtain two solutions, of radii  $0 < \rho_1 < \rho_2$ , say. The right-hand side of (4.3.179) is negative for  $0 < \rho := \sqrt{x^2 + y^2} < \rho_1$ , but positive for  $\rho_1 < \rho < \rho_2$  and negative again beyond  $\rho_2$ . Thus, the orbit at  $\rho_1$  is repelling whereas the one at  $\rho_2$  is attracting. When  $\lambda$  increases to 0, the repelling periodic orbit at  $\rho_1$  moves into the attracting fixed point at 0. When  $\lambda$  then becomes positive, both the repelling periodic orbit and the attracting fixed point disappear, or, more precisely, the latter turns into a repelling fixed point. Only the attracting periodic orbit at  $\rho_2$  remains. The solution of our system of ODEs then has no option but to move away from the no longer attracting fixed point at 0 to the periodic orbit at  $\rho_2$ .

The first bifurcation, the one of (4.3.170), (4.3.171), where a stable fixed point continuously changed into a stable periodic orbit was a so-called supercritical Hopf bifurcation. In contrast to this, in a subcritical Hopf bifurcation, as exemplified by (4.3.177), (4.3.178), an unstable periodic orbit coalesces into a stable



fixed point so that the latter becomes repelling and no stable orbit is present anymore in its vicinity when the relevant parameter passes the bifurcation value. The linearization at 0 is the same for both examples, the supercritical and the subcritical Hopf bifurcation. The linearization possesses a pair of complex conjugate eigenvalues whose real parts vanish at the bifurcation point. In fact, by the theorem of E.Hopf, this is precisely the criterion for such a bifurcation where a stable fixed point bifurcates into a family of periodic orbits.

The preceding examples essentially cover the qualitative types of behaviour of two-dimensional systems of ODEs. This is essentially a consequence of the principle observed as a corollary of the Picard-Lindelöf theorem that two orbits can never intersect or merge. In higher dimensions, however, (even though that principle is still in force) the behaviour can get more complicated, and in fact defies a complete classification.

In any case, however, we have an existence result of the type of Theorem 4.3.2.

**Theorem 4.3.3.** *We consider for  $x = (x^1, \dots, x^n)$  and  $f = (f^1, \dots, f^n)$ ,*

$$\dot{x} = f(x) \text{ for } t \geq 0, \text{ with initial condition } x(0) = x_0 \quad (4.3.180)$$

*and assume that there exist numbers  $m^\alpha < M^\alpha$  for  $\alpha = 1, \dots, n$  with*

$$f(x^1, \dots, x^{\alpha-1}, m^\alpha, x^{\alpha+1}, \dots, x^n) > 0 \text{ and } f(x^1, \dots, x^{\alpha-1}, M^\alpha, x^{\alpha+1}, \dots, x^n) < 0 \quad (4.3.181)$$

*whenever  $m^\beta \leq x^\beta \leq M^\beta$  for all  $\beta = 1, \dots, n$ , and also*

$$m^\beta \leq x_0^\beta \leq M^\beta. \quad (4.3.182)$$

*Then (for a Lipschitz continuous  $f$ ), the solution of (4.3.155) exists for all  $t \geq 0$ .*

The *proof* proceeds as the one of Theorem 4.3.2.

We now exhibit several systems of ODEs that are important as models at various biological scales.

1) Biochemical kinetics:

References here are [28], [24]. The basis here is the law of mass action which states the reaction rate of a chemical reaction is proportional to the concentrations of the reactants raised to the number in which they enter the reaction. That expression is proportional to the collision probability for the reactants. For the simple reaction



when  $k_+$  is the rate constant for the forward reaction that converts  $S_1 + S_2$  into  $2P$  and  $k_-$  is the one for the backward reaction and if we denote the respective concentrations by  $s_1, s_2, p$ , then

$$\dot{s}_1 = \dot{s}_2 = -k_+ s_1 s_2 + k_- p^2 \quad (4.3.184)$$

$$\dot{p} = 2(k_+ s_1 s_2 - k_- p^2). \quad (4.3.185)$$

Enzymatic reactions are of particular importance. The prototype is



Here, the substrate  $S$  and the enzyme  $E$  first form the enzyme-substrate complex  $ES$  in a reversible manner with forward and backward rate constants  $k_1, k_{-1}$ , resp., and then the product  $P$  is irreversibly released from the enzyme  $E$  with rate constant  $k_2$ . When we denote the concentrations of  $E, S, ES, P$  by  $e, s, c, p$ , resp., we obtain the system of ODEs

$$\dot{s} = -k_1es + k_{-1}c \quad (4.3.187)$$

$$\dot{e} = -k_1es + (k_{-1} + k_2)c \quad (4.3.188)$$

$$\dot{c} = k_1es - (k_{-1} + k_2)c \quad (4.3.189)$$

$$\dot{p} = k_2c. \quad (4.3.190)$$

We observe that  $p$  does not appear on the r.h.s of this system. Thus, we need only solve the first 3 equations.  $p$  then is obtained by a simple integration. Moreover, the second and third equations are dependent, and we conclude

$$e(t) + c(t) \equiv e_0 \quad (4.3.191)$$

a constant.

Based on the small amount of enzyme needed for such reactions, the Michaelis-Menten theory makes the assumption of a quasi-steady state for the complex  $ES$ ,

$$\dot{c} = 0. \quad (4.3.192)$$

This is mathematically not unproblematic and requires a singular perturbation analysis, see [28], but here we simply observe the consequence

$$c = \frac{k_1e_0s}{k_1s + k_{-1} + k_2}. \quad (4.3.193)$$

We now move from the molecular to the cellular level and as our next example consider the

2) Hodgkin-Huxley model for the firing of neurons:

The main variable is the potential  $V$  of the neuron, satisfying the ODE

$$C \frac{dV}{dt} = I_e - I_i \quad (4.3.194)$$

where  $C$  is the capacitance of the membrane and  $I_e$  and  $I_i$  are the external and internal currents. The internal current in turn satisfies the equation

$$I_i = g_0(V - V_0) + g_1m^3h(V - V_1) + g_2n^4(V - V_2), \quad (4.3.195)$$

where  $g_0, g_1, g_2 > 0$  and  $V_0, V_1, V_2$  are constants whereas  $m, n, h$  are gating variables corresponding to activation of sodium ( $\text{Na}^+$ ), activation of potassium

(K<sup>+</sup>), and inactivation of Na<sup>+</sup>, resp. Normalizations are such that the gating variables always take their values between 0 and 1 so that they can be interpreted as the probabilities for the corresponding type of channel to be open.

(4.3.194) and (4.3.195) combine to become

$$C \frac{dV}{dt} = I_e - (g_0(V - V_0) + g_1 m^3 h (V - V_1) + g_2 n^4 (V - V_2)). \quad (4.3.196)$$

Whereas  $I_e$  is treated as an external parameter, the internal dynamical regimes crucially depend on the signs of the three terms in (4.3.196) contributing to  $I_i$ . Before going into details, we then formulate a fourth principle for the investigation of systems of ODEs:

4. determine the signs of the diverse summands into which  $f^i$  may be decomposed on the right hand side of (4.3.145) and assess their contribution on the global behaviour of the solution.

Before proceeding, however, we need to clarify the roles of the gating variables.  $m, n, h$  satisfy differential equations of the form

$$\tau_y(V) \frac{dy}{dt} = y_\infty(V) - y \quad (4.3.197)$$

with the limiting value  $y_\infty(V)$  and the time constant  $\tau_y(V)$  indicating the time scale on which the corresponding gating variable varies. A simpler model would consist of taking  $y$  directly as a function of  $V$ , equal to the equilibrium value, that is,  $y = y_\infty(V)$ . The model of Hodgkin-Huxley instead introduces some additional temporal dynamics where  $y$  relaxes to that equilibrium value on the time scale described by  $\tau_y(V)$ . Thus, in particular, it does not follow changes in  $V$  instantaneously, but needs some time to adapt.

The Hodgkin-Huxley system then consists of 4 differential equations, namely (4.3.196) for the voltage  $V$  and three equations of type (4.3.197) for the three gating variables.

It is important for the dynamics of the Hodgkin-Huxley model that while  $m_\infty$  and  $n_\infty$  are increasing functions of  $V$  ( $m_\infty$  starts to rise only at a somewhat higher value of  $V$  (around -80 mV) than  $n_\infty$ ),  $h_\infty$  is a decreasing function. Moreover, the time constant  $\tau_m$  is much smaller than the time constants  $\tau_n, \tau_h$  (which peak at values of  $V$  between -80 and -70 mV), and so  $m$  changes much faster than  $n$  and  $h$ , in fact on the same scale as  $V$ .

The reversal potentials in (4.3.195) are

$$V_1 = 50 \text{mV} \quad (4.3.198)$$

$$V_2 = -77 \text{mV} \quad (4.3.199)$$

$$V_0 = -54.4 \text{mV}. \quad (4.3.200)$$

We now present a qualitative discussion of the dynamics of the Hodgkin-Huxley model. Suppose the system initially is at rest near  $V_0$ . Then  $h = h_\infty(V_0)$  and  $n = n_\infty(V_0)$  are positive (in the order of magnitude 1/2) while  $m = m_\infty(V_0)$  is

close to zero. The relevant term in (4.3.196) then is  $g_0(V - V_0)$  which stabilizes the rest point  $V_0$ . If now some positive current  $I_e$  is injected, a positive feedback dynamics between  $V$  and  $m$  sets in, as in the range we are entering they are each increasing functions of the other one (recall that  $m_\infty$  is an increasing function of  $V$ ). Namely, once  $V$  rises to about -50mV,  $m$  suddenly rises to significantly positive values, and as  $h$  is also positive, the  $\text{Na}^+$  term causes a sharp decrease in the interior current  $I_i$  and thus a further rapid increase in  $V$ , up to the  $\text{Na}^+$  equilibrium value of 50mV. Thus, the potential  $V$  has risen from about -50mV to 50mV within a very short time period. This event is called a spike. However, as  $V$  rises,  $h$  decreases towards 0, and so the  $\text{Na}^+$  current gets deactivated. In that entire sequence, from the initial rise of  $m$  until the decrease of  $h$ , the dynamics is essentially driven by the term  $g_1 m^3 h (V - V_1)$  in (4.3.196). That term also ensures that the voltage does not exceed the peak value  $V_1$ . Moreover,  $n$  increases, and so the  $\text{K}^+$  is activated more strongly, and this causes a decrease of  $V$  even below the resting value  $V_0$ , down to about  $V_2$ , a hyperpolarization. The crucial term for the  $V$  dynamics now is  $g_2 n^4 (V - V_2)$ . This causes a refractory period during which no further spike can be fired, during which (in the absence of a further external current) all variables are readjusted back to their resting values.

It is important to note that already a relatively small or short external current that is barely able to increase  $V$  by about 5mV suffices to trigger the spiking of the neuron, that is, an increase of  $V$  by about 100mV. Thus, a neuron is a device that can amplify the effect of an external input. This input is usually transmitted to a neuron via synaptic connections from other neurons, and one can then study the spreading of activation in a network of neurons.

For more details on the Hodgkin-Huxley model, see [28, 25] and, for new mathematical aspects of it, [32].

The Hodgkin-Huxley model is one of the very few biophysical models that not only captures a qualitative behavior, but allows for numerically accurate predictions. It is somewhat complicated, however, in the sense that it is not easy to assess the effects of variations of the parameters involved and that systems of connected Hodgkin-Huxley type neurons become very difficult to analyze. Therefore, at the expense of numerical accuracy, one may seek a simplified model that still captures the important qualitative aspects of spiking neurons. Thus, one seeks a simpler system with the same qualitative behavior of its solutions as the Hodgkin-Huxley model. This starts from the observation that the 4 dependent variables of the Hodgkin-Huxley system evolve on 2 different time scales, a fast one for the evolution of  $V$  and  $m$  (which both return rapidly to their rest states after a spike) and a slower one for  $n$  and  $h$ . In particular, since  $m$  changes on the same time scale as  $V$  itself, it can be taken as a function of the latter and essentially be eliminated from the system. Therefore, one lumps  $V$  and  $m$  together as a single variable  $v$ , and  $n$  and  $1 - h$  (which show similar behavior) as  $w$ . This leads to the FitzHugh-Nagumo system<sup>5</sup> (where we

---

<sup>5</sup>Here, we do not give a detailed derivation of the FitzHugh-Nagumo system from the

abbreviate  $\dot{v} = \frac{dv}{dt}$  etc.)

$$\dot{v} = v(a - v)(v - 1) - w + \lambda \quad (4.3.201)$$

$$\dot{w} = bv - cw \quad (4.3.202)$$

with constants  $0 < a < 1$ ,  $b, c > 0$ . The parameter  $\lambda$  here represents the external current  $I_e$ , i.e. the input to the neuron.

Thus, the term with  $m^3$  in (4.3.195) translates into the cubic term in (4.3.201). Since the leading coefficient of that cubic term is negative, the dynamics is always confined to some bounded region in the  $v, w$ -plane.  $w$  enters the system only linearly, and in turn its own evolution equation is linear. The rest point  $V = V_0$  in (4.3.195) becomes the origin in (4.3.201), (4.3.202).

For this system, the qualitative aspects of the dynamics can be readily analyzed (see e.g. [28]). We abbreviate

$$f(v) := v(a - v)(v - 1). \quad (4.3.203)$$

Following the general strategy outlined above for the qualitative analysis of a system of ODEs, we identify the rest points; this is achieved by putting all the time derivatives, that is, the left hand sides of (4.3.201), (4.3.202) equal to 0 and solving the resulting algebraic equation. We start with the analysis for  $\lambda = 0$ . In that case, the rest points for the FitzHugh-Nagumo system are determined by the equations

$$0 = f(v) - w \quad (4.3.204)$$

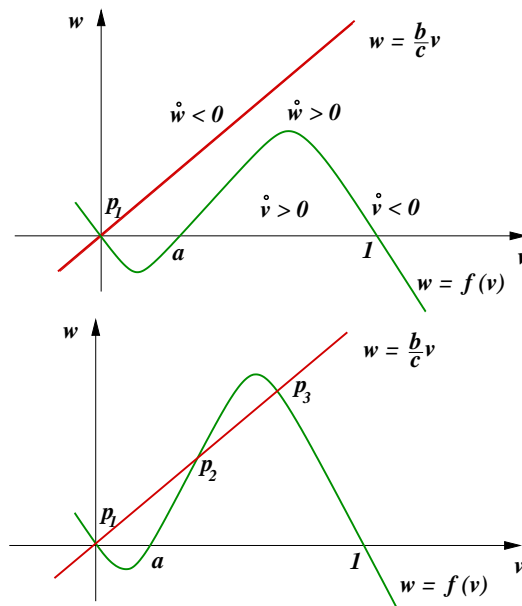
$$0 = bv - cw. \quad (4.3.205)$$

Depending on the values of the parameters  $a, b, c$ , the behavior is described by one of the following two figures

Thus, the rest points are at the intersections of the red ( $w = \frac{b}{c}v$ ) and the green ( $w = f(v)$ ) curves. For generic parameter values, we either find one stable fixed point  $P_1$ , or two stable ones  $P_1, P_3$  and one unstable one  $P_2$ . In the vicinity of the stable fixed point, the behavior of the system is determined by the quadratic term  $av^2$  of  $f(v)$  while the negative cubic term  $-v^3$  becomes effective only for large values, leading to a resetting of  $v$ . Thus, small perturbations of  $P_1$  asymptotically return to  $P_1$ . If at  $w = 0$ , however,  $v$  is thrown above the value  $a$ , one gets into the region  $\dot{v} > 0$ , and  $v$  thus increases until returning again to the green curve. As one also is in the region  $\dot{w} > 0$ ,  $w$  increases as well until reaching the red curve. In the situation captured in the second figure, one then approaches the second stable fixed point  $P_3$ . In the situation of the first figure, however, one gets into the region where  $\dot{v}$  and  $\dot{w}$  both are negative, and  $v$  and  $w$  thus decrease. The dynamics then gets into the region  $v < 0$  left and above the red curve, until  $\dot{v}$  eventually becomes positive again, and the dynamics returns to the starting point  $P_1$ . This process then is interpreted as the firing of the neuron when the threshold  $a$  is exceeded. In summary, we see a

---

Hodgkin-Huxley one. See e.g. [28, 25].



qualitatively different behavior, depending on whether the initial perturbation is small or large. In the first case,  $v$  directly decreases to its rest point. In the second case, it needs to increase first above a certain value before it is able to return to the rest point.

We now wish to analyze the role of the parameter  $\lambda$  that has been left out of the picture so far. After introducing  $\lambda$ , the curve  $\dot{v} = 0$  becomes

$$w = f(v) + \lambda; \quad (4.3.206)$$

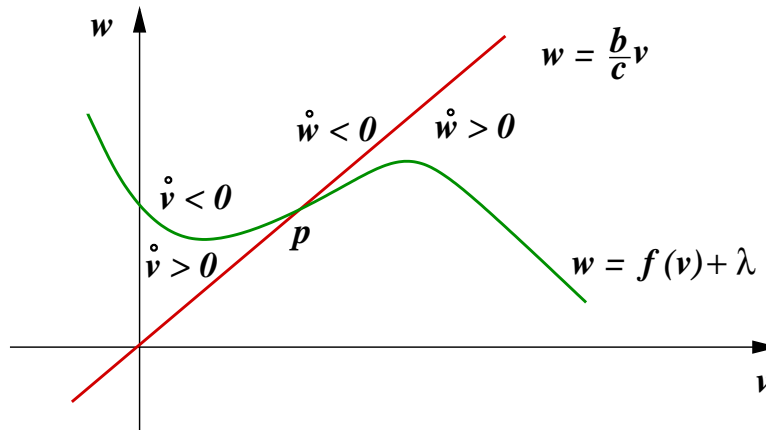
thus, the green curve is shifted. If  $\lambda$  is positive, from the situation of the first figure, we can either get into the one of the second figure (a transition representing a *saddle-node bifurcation* as a stable and an unstable rest point emerge from a contact point between the two curves), or into the one of the figure below

ACHTUNG: BILDER MÜSSEN NOCH DURCH TYPISCHE TRAJEKTORIEN ERWEITERT WERDEN

In that scenario, the single fixed point  $P$  is unstable, and perturbations from the rest position lead first away from  $P$  and then turn into oscillations around that rest position as the asymptotic behavior is dominated by the cubic term. We see a Hopf bifurcation, as described above.

Our next example is relevant for a much larger scale, the one of ecological interactions of populations:

3) The Lotka-Volterra system for the sizes  $x^i$  of  $n$  interacting populations (good



references being [19, 20]) is

$$\dot{x}^i = x^i (a_i + \sum_{j=1}^d b_{ij} x^j) \text{ for } i = 1, \dots, n. \quad (4.3.207)$$

$a_i$  is intrinsic growth or decay rate of the  $i$ th population in the absence of the other populations, and  $b_{ij}$  is the strength of the effect that the  $j^{th}$  population has on the  $i^{th}$  one.  $a_1$  is positive (negative) iff  $x^i$  has an inherent tendency to grow (decay), and  $b_{ij}$  is positive (negative) iff  $x^j$  enhances (inhibits) the growth of  $x^i$ , e.g. if population  $i$  feeds on (is preyed upon by) population  $j$ ; both  $b_{ij}$  and  $b_{ji}$  are negative if the two corresponding populations compete. The self-effect  $b_{ii}$  is typically negative, expressing a limiting carrying capacity of the environment or interspecific competition for resources, or at least non-positive. Thus, when  $x^i$  gets too large, this term takes over and keeps the population in check. Biological and other populations always satisfy

$$x^i(t) \geq 0. \quad (4.3.208)$$

Thus, we only need to investigate solutions in the positive quadrant.

For a single population, we consider the logistic or Fisher equation (see (4.3.151))

$$\dot{x}(t) = x(a + bx) \text{ with } a > 0, b < 0. \quad (4.3.209)$$

This is about a population growing under the condition of limited or constrained resources, so that, when it gets too large, the capacity limits take over and keep it in balance.  $x = 0$  is an unstable fixed point,  $x = -a/b$  a stable one.

For the case of two populations, there are three non-trivial scenarios:

1. Predator-prey or parasitism: Population 1 is the prey or host, population 2 its predator or parasite:

$b_{12} < 0$  (the prey is fed upon by the predators)

$b_{21} > 0$  (the presence of prey leads to growth of the predator population).

In the predator-prey scenario, one also typically has

- $a_1 > 0$  (the prey population grows in the absence of predators)
- $a_2 < 0$  (the predator population decays in the absence of prey).

2. Competition

- $b_{12} < 0$  and  $b_{21} < 0$  (the two populations inhibit each other).

3. Symbiosis:

- $b_{12} > 0$  and  $b_{21} > 0$  (the two populations support each other).

We now look at the example of the two-dimensional predator-prey model without intraspecific competition, that is, we have a prey population of size  $x^1$  and a predator population of size  $x^2$ ,

$$\begin{aligned}\dot{x}^1 &= x^1(a_1 + b_{12}x^2) \\ \dot{x}^2 &= x^2(a_2 + b_{21}x^1),\end{aligned}\tag{4.3.210}$$

with

$$a_1 > 0, \quad a_2 < 0, \quad b_{12} < 0, \quad b_{21} > 0.$$

We first observe that  $(x^1, x^2) = (0, 0)$  is a fixed point. Linearization shows that this fixed point is a saddle. On the  $x^1$ -axis, the solution expands according to  $x^1(t) = x^1(0)e^{a_1 t}$ ,  $x^2(t) = 0$ , whereas it contracts along the  $x^2$ -axis as  $a_2 < 0$ ,  $x^1(t) = 0$ ,  $x^2(t) = x^2(0)e^{a_2 t}$ . In particular, since the two axes are orbits, the solution cannot cross them, that is, when starting with non-negative values, it will never turn negative, in accordance with (4.3.208).

Another fixed point is

$$\bar{x}^1 = -\frac{a_2}{b_{21}}, \quad \bar{x}^2 = -\frac{a_1}{b_{12}}\tag{4.3.211}$$

All the other orbits in the positive quadrant are periodic, circling this fixed point counterclockwise. This is seen either from the local behavior of the trajectories near  $(0, 0)$  or by looking at

$$V(x^1, x^2) := b_{21}(\bar{x}^1 \log x^1 - x^1) - b_{12}(\bar{x}^2 \log x^2 - x^2),\tag{4.3.212}$$

which satisfies

$$\begin{aligned}\frac{d}{dt}V(x^1(t), x^2(t)) &= -a_2 \frac{\dot{x}^1}{x^1} - b_{21}\dot{x}^1 + a_1 \frac{\dot{x}^2}{x^2} + b_{12}\dot{x}^2 \text{ by (4.3.211)} \\ &= 0 \text{ by (4.3.212)}.\end{aligned}$$

Thus,  $V(x^1, x^2)$  is a constant of motion.  $V$  attains its unique maximum at  $(\bar{x}^1, \bar{x}^2)$ , and so the curves  $V(x^1, x^2) \equiv \text{constant}$  are circles, that is, closed curves, around this point. The motion on such a circle is counterclockwise because in



the case  $x^1(t) > \bar{x}^1$ ,  $x^2(t) > \bar{x}^2$  for example, we have  $\dot{x}^1(t) < 0$ ,  $\dot{x}^2(t) > 0$ . On the line  $x^1 = \bar{x}^1$ ,  $\dot{x}^2(t) = 0$ , that is,  $x^2$  stays constant there, and on  $x^2 = \bar{x}^2$ ,  $x^1$  stays constant.

Thus, the prey and predator populations oscillate periodically in this model.

The behaviour of the preceding system with its family of periodic orbits is not stable under small perturbations, as we already know from our discussion of the Hopf bifurcation above. For example when we include intraspecific competition, obtaining the system

$$\begin{aligned}\dot{x}^1 &= x^1(a_1 + b_{11}x^1 + b_{12}x^2) \\ \dot{x}^2 &= x^2(a_2 + b_{21}x^1 + b_{22}x^2)\end{aligned}\quad (4.3.213)$$

where we now assume

$b_{11} < 0$  (the members of the prey population compete for food or other resources)  
 $b_{22} \leq 0$ ,

the qualitative behaviour of the system becomes different. We find a second fixed point on the positive  $x^1$ -axis,  $(-\frac{a_1}{b_{11}}, 0)$ . This fixed point is always attractive for  $x^1$ , because in case  $x^2(t) = 0$ , we have the logistic equation (4.3.209)

$$\dot{x}^1(t) = x^1(a_1 + b_{11}x^1) \text{ with } a_1 > 0, b_{11} < 0. \quad (4.3.214)$$

It is also attractive for  $x^2$  when

$$a_2b_{11} - a_1b_{21} > 0.$$

In that case, there is no other fixed point in the positive quadrant, and in fact for any solution

$$\lim_{t \rightarrow \infty} x^2(t) = 0.$$

The predator becomes extinct. Thus, we conclude that a small intraspecific competition among the prey population may lead to the extinction of their predators.

If, however,

$$a_2b_{11} - a_1b_{21} < 0,$$

then

$$\begin{aligned}\bar{x}^1 &= \frac{a_2b_{12} - a_1b_{22}}{b_{11}b_{22} - b_{12}b_{21}} > 0 \\ \bar{x}^2 &= \frac{a_1b_{21} - a_2b_{11}}{b_{11}b_{22} - b_{12}b_{21}} > 0\end{aligned}$$

is a fixed point in the positive quadrant.

With  $V(x^1, x^2)$  as in (4.3.212),

$$\frac{d}{dt}V(x^1(t), x^2(t)) = -b_{11}b_{21}(\bar{x}^1 - x^1(t))^2 + b_{12}b_{22}(\bar{x}^2 - x^2(t))^2 > 0,$$

unless  $(x^1, x^2) = (\bar{x}^1, \bar{x}^2)$  in which case this derivative vanishes. Thus,  $V(x^1(t), x^2(t))$  increases along every orbit, and equilibrium is possible only at its maximum, at the fixed point  $(\bar{x}^1, \bar{x}^2)$ . The orbits in the positive quadrant then all spiral counterclockwise towards this fixed point. Thus, in this case, the two populations eventually converge to this fixed point.

We have observed that an arbitrarily small variation of the original system, here by introducing competition among the hosts, changes the global qualitative behavior of the solutions. Therefore, one cannot expect that this model leads to a qualitatively robust and structurally stable behavior, and predictions based on such a model need to be examined with great care. Volterra originally introduced the model to explain the periodic oscillations in two fish populations in the Adriatic, one preying upon the other one. As it turned out, however, this periodic behavior is not caused by an interaction of the two populations according to the model, but rather by periodic changes in the water temperature, that is, by external periodic forcing. Although the model therefore fails its original purpose, it has become useful at a more abstract level, for game theoretic models of interactions inside populations, see [19, 20].

## 4.4 Reaction-diffusion systems

References: [34, 28, 21].

### 4.4.1 Reaction-diffusion equations

Let  $\Omega \subset \mathbb{R}^d$  be open and bounded.

We consider the equation

$$\begin{aligned} u_t(x, t) &= \Delta u(x, t) + f(x, t, u) \text{ for } x \in \Omega, 0 < t < T & (4.4.1) \\ u(x, 0) &= \phi(x) \text{ for } x \in \Omega \\ u(y, t) &= g(y, t) \text{ for } y \in \partial\Omega, 0 < t < T \end{aligned}$$

for continuous and bounded initial and boundary values and a differentiable reaction term  $f$ .

We can consider this is a generalization of

1. either the ODE

$$u_t(t) = f(t, u), \quad (4.4.2)$$

(at least in the case where the function  $f$  in (4.4.1) does not depend on  $x$ ), that is, of an equation that does not depend on the spatial variable  $x$  and therefore describes a spatially homogeneous state,

2. or the linear heat equation

$$u_t(x, t) = \Delta u(x, t), \quad (4.4.3)$$

that is, of an equation that describes a linear diffusion process in space.

The first equation describes a reaction process, and thus, a reaction-diffusion equation models reaction processes taking place at all points in space simultaneously and being diffusively coupled. It turns out that such an interplay of reaction and diffusion processes can lead to more interesting patterns than either of these processes alone.

In order to understand the relationship between the two processes better, we can also introduce a diffusion coefficient  $d$  and consider the more general equation

$$u_t(x, t) = d\Delta u(x, t) + f(x, t, u), \quad u(x, 0) = \phi(x) \quad (4.4.4)$$

When  $d = 0$ , we have a system of ODEs indexed by the points  $x$ , but without any coupling or interaction between those points. Thus, at each point  $x$ , the dynamics is driven by the reaction term, and the result depends only on the initial condition  $\phi(x)$  at that particular  $x$ . When we let  $d \rightarrow \infty$ , we obtain an equation for the spatially integrated variables, that is, a single ODE for the spatially averaged quantity, and there will be no variation between the different points  $x$ . In general, in physical and biological processes, conservation rules will limit growth, and growth at one point  $x$  then is only possible at the expense of other points. In principle, a reaction-diffusion equation can then lead to optimal resource allocation in the limit  $t \rightarrow \infty$ . The time scale on which this takes place will depend on the diffusion coefficient  $d$ .

Recalling (4.3.155), let us consider

$$u_t = \Delta u + u(1 - u), \quad (4.4.5)$$

for  $u = u(x, t)$ , that is, the logistic (Verhulst, Fisher) equation with a diffusion term. This equation is sometimes called the Kolmogorov-Fisher equation. It can serve as a model for a population in a uniform habitat with limited capacity that reproduces and diffuses in space. As such, we expect that the behavior of a solution  $u(x, t)$  at a point  $x$  is not very different from the solution of the ODE  $y_t(t) = y(t)(1 - y(t))$  with initial value  $y(0) = u(x, 0)$ . In particular, we expect that a solution  $u(x, t)$  with a positive initial value converges to 1 for  $t \rightarrow \infty$ , unless this is prevented by the boundary condition. However, there will be diffusion between the different  $x \in \Omega$ , and this may decrease the differences in their respective initial conditions faster than the dynamical evolution by the reaction term alone. This effect becomes more important when we consider the spatially inhomogeneous equation

$$u_t(x, t) = \Delta u(x, t) + u(x, t)(a(x) - b(x)u(x, t)) \quad (4.4.6)$$

for positive functions  $a, b$ . Thus, the intrinsic growth rate and the capacity limitations depend on  $x$ . Now, the stable equilibrium point for the reaction term,  $u = \frac{a(x)}{b(x)}$ , depends on the spatial position  $x$ . Without diffusion, at every  $x$  then in the limit  $t \rightarrow \infty$ , this equilibrium would be obtained. With diffusion, however, we expect some harmonization between the higher and lower values of that equilibrium.

In applications, the dependent variable  $u$  typically describes some density, and therefore only non-negative solutions  $u$  will be meaningful.

**Example:**

$$\begin{aligned} u_t(x, t) &= \Delta u(x, t) + u^2(x, t) \text{ for } x \in \Omega, 0 < t & (4.4.7) \\ u(x, 0) &= \phi(x) \text{ for } x \in \Omega \\ u(y, t) &= 0 \text{ for } y \in \partial\Omega, 0 < t \end{aligned}$$

with

$$\phi > 0 \text{ in } \Omega. \quad (4.4.8)$$

We recall that the ODE

$$u_t(t) = u^2(t) \quad (4.4.9)$$

for positive initial value  $u(0)$  did blow up in finite time (cf. (4.3.153), (4.3.154)), that is, did not possess a solution that exists for all  $t > 0$ . We shall now show that the same happens for (4.4.7) provided the initial values  $\phi$  are sufficiently large. To make that condition precise, we recall the first Dirichlet eigenvalue  $\lambda_1$  and the corresponding eigenfunction  $u_1$  from Theorem 4.1.1, solving

$$\Delta u_1 + \lambda_1 u_1 = 0, \quad u_1 = 0 \text{ on } \partial\Omega,$$

and recall that, by (4.1.69),

$$u_1(x) > 0 \text{ for } x \in \Omega; \quad (4.4.10)$$

we may normalize  $u_1$  to satisfy

$$\int_{\Omega} u_1(x) dx = 1. \quad (4.4.11)$$

By the maximum principle Lemma 4.1.2, because of (4.4.8), we have

$$u(x, t) > 0 \text{ for } x \in \Omega, 0 < t. \quad (4.4.12)$$

We look at the auxiliary function

$$y(t) := \int_{\Omega} u(x, t) u_1(x) dx \quad (4.4.13)$$

which satisfies

$$\begin{aligned} \dot{y}(t) &= \int_{\Omega} u_t(x, t) u_1(x) dx = \int_{\Omega} (\Delta u(x, t) + u^2(x, t)) u_1(x) dx \\ &= \int_{\Omega} (\Delta u_1(x) u(x, t) + u^2(x, t) u_1(x)) dx = -\lambda_1 y(t) + \int_{\Omega} u^2(x, t) u_1(x) dx \\ &\geq -\lambda_1 y(t) + y^2(t) \end{aligned} \quad (4.4.14)$$

since, by the definition (4.4.13), Hölder's inequality,<sup>6</sup> and (4.4.11)

$$y^2(t) \leq \left( \int u^2(x, t)u_1(x)dx \right) \left( \int u_1(x)dx \right) = \int u^2(x, t)u_1(x)dx.$$

When now

$$y(0) = \int_{\Omega} u(x, 0)u_1(x)dx > \lambda_1, \tag{4.4.15}$$

then (4.4.14) easily implies that  $y(t)$  will blow up in finite time, similarly to a solution of (4.4.9) with positive  $x(0)$ . ((In fact, (4.4.15) implies that  $\dot{y}(0) > 0$ , and then subsequently  $\dot{y}(t) > 0$  for all  $t \geq 0$ , and the solution will grow and the quadratic term  $y^2$  will dominate the behavior.) This implies that, when (4.4.15) holds, (4.4.7) cannot possess a smooth solution for all positive  $t$ .

The maximum principle will allow for a comparison of solutions of a reaction diffusion equation:

**Lemma 4.4.1.** *Let  $u, v$  be of class  $C^2$  for  $x \in \Omega$ ,  $0 < t < T$ , and bounded in  $\bar{\Omega} \times [0, T]$ , and satisfy*

$$u_t - \Delta u - f(x, t, u) \geq v_t - \Delta v - f(x, t, v) \text{ for } x \in \Omega, 0 < t < T \tag{4.4.16}$$

$$u(x, 0) \geq v(x, 0) \text{ for } x \in \Omega, \tag{4.4.17}$$

$$u(y, t) \geq v(y, t) \text{ for } y \in \partial\Omega, 0 < t < T, \text{ or alternatively} \tag{4.4.18}$$

$$\frac{\partial u(y, t)}{\partial n} \geq \frac{\partial v(y, t)}{\partial n} \text{ for } y \in \partial\Omega, 0 < t < T. \tag{4.4.19}$$

Then

$$u(x, t) \geq v(x, t) \text{ for } x \in \Omega, 0 \leq t \leq T. \tag{4.4.20}$$

*Proof.*

$$w(x, t) := u(x, t) - v(x, t)$$

is non-negative for  $x \in \partial\Omega, 0 < t < T$  and  $x \in \Omega, t = 0$  and satisfies

$$w_t - \Delta w - f_{\eta}(x, t, \eta)w \geq 0 \tag{4.4.21}$$

for some intermediate  $\eta = \epsilon u + (1 - \epsilon)v, 0 \leq \epsilon \leq 1$ . The function

$$z(x, t) := w(x, t)e^{-\mu t} \tag{4.4.22}$$

then satisfies

$$e^{\mu t}(z_t - \Delta z - f_{\eta}(x, t, \eta)z - \mu z) = w_t - \Delta w - f_{\eta}(x, t, \eta)w \geq 0 \tag{4.4.23}$$

<sup>6</sup>Hölder's inequality says that for two  $L^2$  functions  $f, g$  (that is, functions with finite  $\int f^2(x)dx, \int g^2(x)dx$ ), we have

$$\left( \int f(x)g(x)dx \right)^2 \leq \int f^2(x)dx \int g^2(x)dx.$$

There are also other such calculus inequalities, like the Poincaré and Sobolev ones, which are very useful for the control of solutions of PDEs.

and by making  $\mu$  sufficiently large, since  $f_\eta(x, t, \eta)$  is bounded by the boundedness assumption on  $u, v$ , therefore

$$z_t - \Delta z \geq 0 \quad (4.4.24)$$

as long as  $z \geq 0$ . The strong maximum principle (Lemma 4.1.2) then implies that when  $z$  is non-negative and positive somewhere, it cannot become 0 for some  $x \in \Omega, t > 0$ . Since non-negativity of  $z$  implies non-negativity of the difference  $w$  of our solutions, this yields the claim.  $\square$

The version of this lemma where an inequality between the exterior normal derivatives of the solutions is assumed allows for a comparison of a solution of (4.4.1) where  $f$  does not depend on  $x$ , that is,

$$u_t = \Delta u + f(t, u) \quad (4.4.25)$$

$$u(x, 0) = \phi(x) \quad (4.4.26)$$

$$\frac{\partial u(y, t)}{\partial n} = 0 \text{ on } \partial\Omega \quad (4.4.27)$$

with one of the corresponding ODE

$$y_t = f(t, y) \quad (4.4.28)$$

$$y(0) = y_0. \quad (4.4.29)$$

When, for example,

$$y_0 \leq \phi(x) \text{ for all } x \in \Omega, \quad (4.4.30)$$

then we conclude that

$$y(t) \leq u(x, t) \text{ for all } x \in \Omega. \quad (4.4.31)$$

Similarly, when  $y_0$  is bigger than  $\phi$  in  $\Omega$ , then the corresponding solution  $y$  of (4.4.28) controls the solution  $u$  of (4.4.25) from above.

For example, when

$$f(t, u) = -u^3 \quad (4.4.32)$$

then any solution  $y(t)$  of (4.4.28) goes to zero for  $t \rightarrow \infty$ , and when we can then sandwich any solution  $u(x, t)$  of (4.4.25) between solutions of (4.4.28) with smaller and larger initial values, resp., than the initial values  $\phi$  of  $u$ , and therefore conclude that such a solution  $u$  – if it exists for all time – also tends to 0 for  $t \rightarrow \infty$ .

In fact, the general theory of parabolic equations tells us that a solution exists for all time when it can be shown to be bounded. The latter is precisely what can be achieved by such comparison arguments based on the maximum principle. In particular, when solutions of the corresponding reaction ODE stay bounded – and hence exist for all time – for a range of initial values, then so do solutions of the reaction diffusion equation for the same range of initial values (at least in the case of Neumann boundary conditions, but suitable results also hold for Dirichlet boundary conditions).

### 4.4.2 Travelling waves

We consider the reaction-diffusion equation in one-dimensional space

$$u_t = u_{xx} + f(u) \quad (4.4.33)$$

and look for solutions of the form

$$u(x, t) = v(x - ct) = v(s), \text{ with } s := x - ct. \quad (4.4.34)$$

This travelling wave solution moves at constant speed  $c$ , assumed to be  $> 0$  w.l.o.g, in the increasing  $x$ -direction. In particular, if we move the coordinate system with speed  $c$ , that is, keep  $x - ct$  constant, then the solution also stays constant. We do not expect such a solution for every wave speed  $c$ , but at most for particular values that then need to be determined.

A travelling wave solution  $v(s)$  of (4.4.33) satisfies the ODE

$$v''(s) + cv'(s) + f(v) = 0, \text{ with } ' = \frac{d}{ds}. \quad (4.4.35)$$

When  $f \equiv 0$ , then a solution must be of the form  $v(s) = c_0 + c_1 e^{-cs}$  and therefore becomes unbounded for  $s \rightarrow -\infty$ , that is for  $t \rightarrow \infty$ . In other words, for the heat equation, there is no non-trivial bounded travelling wave. In contrast to this, depending on the precise non-linear structure of  $f$ , such travelling waves solutions may exist for reaction-diffusion equations. This is one of the reasons why such equations are interesting.

**Example:** The Kolmogorov-Fisher equation (4.4.5)

$$u_t = \Delta u + u(1 - u). \quad (4.4.36)$$

It models the spatial spread of a population that grows in an environment with limited carrying capacity. Fisher used it as a model for the spread of an advantageous gene in a population. One then is only interested in non-negative solutions  $u$ .

Here, we consider the case where space has only one dimension,

$$u_t = u_{xx} + u(1 - u). \quad (4.4.37)$$

The fixed points of the underlying reaction equation

$$u_t = u(1 - u) \quad (4.4.38)$$

are  $u = 0$  and  $u = 1$ . The first one is unstable, the second one stable. The travelling wave equation (4.4.35) then is

$$v''(s) + cv'(s) + v(1 - v) = 0. \quad (4.4.39)$$

With  $w := v'$ , this is converted into the first order system

$$v' = w, \quad w' = -cw - v(1 - v). \quad (4.4.40)$$

The fixed points then are  $(0, 0)$  and  $(1, 0)$ . The eigenvalues of the linearization at  $(0, 0)$  (cf. 4.3.1, in particular, the discussion following (4.3.163)) are

$$\lambda_{\pm} = \frac{1}{2}(-c \pm \sqrt{c^2 - 4}). \quad (4.4.41)$$

For  $c^2 \geq 4$ , they are both real and negative, and so we obtain a stable node. For  $c^2 < 4$ , they are conjugate complex with a negative real part, and we obtain a stable spiral. Since a stable spiral oscillates about 0, in that case, we cannot expect a non-negative solution, and so, we do not consider this case here. Also, for symmetry reasons, we may restrict ourselves to the case  $c > 0$ , and since we want to exclude the spiral then to  $c \geq 2$ .

The eigenvalues of the linearization at  $(1, 0)$  are

$$\lambda_{\pm} = \frac{1}{2}(-c \pm \sqrt{c^2 + 4}); \quad (4.4.42)$$

they are real and of different signs, and we obtain a saddle. Thus, the stability properties are reversed when compared to (4.4.38) which, of course, results from the fact that  $\frac{ds}{dt} = -c$  is negative.

For  $c \geq 2$ , one finds a solution with  $v \geq 0$  from  $(1, 0)$  to  $(0, 0)$ , that is, with  $v(-\infty) = 1, v(\infty) = 0$ .  $v' \leq 0$  for this solution. We recall that the value of a travelling wave solution is constant when  $x - ct$  is constant. Thus, in the present case, when time  $t$  advances, the values for large negative values of  $x$  which are close to 1 are propagated to the whole real line, and for  $t \rightarrow \infty$ , the solution becomes 1 everywhere. In this sense, the behavior of the ODE (4.4.38) where a trajectory goes from the unstable fixed point 0 to the stable fixed point 1 is translated into a travelling wave that spreads a nucleus taking the value 1 for  $x = -\infty$  to the entire space.

The question for which initial conditions a solution of (4.4.37) evolves to such a travelling wave, and what the value of  $c$  then is, has been widely studied in the literature since the seminal work of Kolmogorov and his coworkers [26]. For example, they showed when  $u(x, 0) = 1$  for  $x \leq x_1$ ,  $0 \leq u(x, 0) \leq 1$  for  $x_1 \leq x \leq x_2$ ,  $u(x, 0) = 0$  for  $x \geq x_2$ , then the solution  $u(x, t)$  evolves towards a travelling wave with speed  $c = 2$ . In general, the wave speed  $c$  depends on the asymptotic behavior of  $u(x, 0)$  for  $x \rightarrow \pm\infty$ . Under the assumptions just mentioned, the solution thus converges to the travelling wave with minimal wave speed. To make this more precise, we consider the simplest case

$$u(x, 0) = 1 \text{ for } x < 0, \quad u(x, 0) = 0 \text{ for } x \geq 0. \quad (4.4.43)$$

For each  $t > 0$ , we then find a unique  $\theta(t)$  by

$$u(\theta(t), t) = \frac{1}{2}. \quad (4.4.44)$$

$u(x, t)$  then converges to the travelling wave  $v(x - 2t)$  in the sense that

$$u(x + \theta(t), t) \rightarrow v(x) \text{ for all } x \in \mathbb{R} \text{ as } t \rightarrow \infty. \quad (4.4.45)$$



Since  $u_x \theta'(t) + u_t = 0$  and we have the asymptotic relationship  $\frac{u_t}{u_x} \rightarrow \frac{-2v'}{v'} = -2$  for our wave speed  $c = 2$ , we also obtain

$$\theta'(t) \rightarrow 2 \text{ for } t \rightarrow \infty. \quad (4.4.46)$$

This also yields some insights for the asymptotic analysis when we let the diffusion speed go to 0, that is, consider instead of (4.4.37)

$$u_t = \epsilon u_{xx} + u(1 - u). \quad (4.4.47)$$

In this case, assuming again (4.4.43), the minimal wave speed is  $c = 2\sqrt{\epsilon}$ . The solution  $u$  then converges to a travelling wave  $v_\epsilon(x - 2\sqrt{\epsilon}t)$  with that wave speed, that is,  $u(x + \theta(t), t) \rightarrow v_\epsilon(x)$  and  $\theta'(t) \rightarrow 2\sqrt{\epsilon}$  for  $t \rightarrow \infty$ . Thus, of course, the wave front moves slower and slower as the diffusion speed decreases.

Alternatively, one may look at  $y(x, \tau) := u(x, \frac{\tau}{\epsilon})$  which then solves

$$y_\tau = y_{xx} + \frac{1}{\epsilon} y(1 - y) \quad (4.4.48)$$

leading to the minimal wave speed  $\frac{2}{\sqrt{\epsilon}}$  with which the solution asymptotically moves. Here, we have rescaled the time so that the diffusion speed stays the same, but the reaction term then explodes in the limit, causing faster and faster wave front propagation.

### 4.4.3 Reaction-diffusion systems

Different diffusion coefficients (one coefficient  $\ll 1$ , another  $\gg 1$ ): one variable can build spatial concentrations and adapts slowly whereas the other, fast adapting one gets enslaved.

We now consider systems of reaction-diffusion equations that are coupled through the reaction terms. These are systems of the form

$$u_t^\alpha(x, t) - d_\alpha \Delta u^\alpha(x, t) = F^\alpha(x, t, u) \quad \text{for } x \in \Omega, t > 0, \alpha = 1, \dots, n. \quad (4.4.49)$$

Here,  $u = (u^1, \dots, u^n)$  has  $n$  components, and the diffusion coefficients  $d_\alpha$  are non-negative constants. When some  $d_\alpha = 0$ , the corresponding equation reduces to an ordinary differential equation for  $u^\alpha$  as a function of time  $t$ ; through the reaction term, it will still be coupled to the other components of  $u$ , however, which satisfy partial differential equations when their diffusion coefficients are positive.

We now state precise existence theorems for the initial-boundary value problem for (4.4.49). The first one is concerned with the existence of solutions for a short time. It is the analogue of the Picard-Lindelöf Theorem 4.3.1. The idea of the proof again is some construction principle. In addition, the maximum principle (see Lemma 4.1.2) plays an important role. This is the reason why now, in contrast to Theorem 4.3.1, we get existence only in forward time.

**Theorem 4.4.1.** *Let, as always, the diffusion constants  $d_\alpha$  all be nonnegative. Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain of class  $C^2$ , and let*

$$g \in C^0(\partial\Omega \times [0, t_0], \mathbb{R}^n), \quad f \in C^0(\bar{\Omega}, \mathbb{R}^n),$$

$$\text{with } g(x, 0) = f(x) \quad \text{for } x \in \partial\Omega,$$

and let

$$F \in C^0(\bar{\Omega} \times [0, t_0] \times \mathbb{R}^n)$$

be Lipschitz continuous w.r.t.  $u$ , that is, there exists a constant  $L$  with

$$|F(x, t, u_1(x)) - F(x, t, u_2(x))| \leq L|u_1(x) - u_2(x)| \quad (4.4.50)$$

for  $x \in \bar{\Omega}$ ,  $t \in [0, t_0]$ ,  $u_1, u_2 \in \mathbb{R}^n$ .

Then there exists some  $0 < t_1 \leq t_0$  for which the initial boundary value problem

$$u_t^\alpha(x, t) - d_\alpha \Delta u^\alpha(x, t) = F^\alpha(x, t, u) \quad \text{for } x \in \Omega, 0 < t \leq t_1, \alpha = 1, \dots, n, \quad (4.4.51)$$

$$\begin{aligned} u(x, t) &= g(x, t) && \text{for } x \in \partial\Omega, 0 < t \leq t_1, \\ u(x, 0) &= f(x) && \text{for } x \in \Omega, \end{aligned} \quad (4.4.52)$$

admits a unique solution that is continuous on  $\bar{\Omega} \times [0, t_1]$ .

The next result addresses the issue of long-time existence.

**Theorem 4.4.2.** *We assume that the preceding assumptions hold for all  $t_0 < \infty$ . We assume furthermore that the following a-priori bound holds: the solution  $u(x, t) = (u^1(x, t), \dots, u^n(x, t))$  of (4.4.51) satisfies the a-priori bound*

$$\sup_{x \in \bar{\Omega}, 0 \leq \tau \leq t} |u(x, \tau)| \leq K \quad (4.4.53)$$

for all times  $t$  for which it exists, with some fixed constant  $K$ . Then the solution  $u(x, t)$  exists for all times  $0 \leq t < \infty$ .

This naturally leads to the question under which assumptions such an a-priori bound holds. This is answered by the analogue of Theorem 4.3.3.

**Theorem 4.4.3.** *Under the above assumptions, suppose that the initial values  $f$  and the boundary values  $g$  both satisfy*

$$m_\alpha \leq f^\alpha(x), g^\alpha(x, t) \leq M_\alpha \quad (4.4.54)$$

for all  $x, t$  where the numbers  $m_\alpha, M_\alpha$  satisfy

$$F^\alpha(x, t, u^1, \dots, u^{\alpha-1}, m_\alpha, u^{\alpha+1}, \dots, u^n) > 0 \quad (4.4.55)$$

$$F^\alpha(x, t, u^1, \dots, u^{\alpha-1}, M_\alpha, u^{\alpha+1}, \dots, u^n) < 0 \quad (4.4.56)$$

whenever  $m_\alpha \leq u^\alpha \leq M_\alpha$  for  $\alpha = 1, \dots, n$ . Then we have the a-priori bounds

$$m_\alpha \leq u^\alpha(x, t) \leq M_\alpha \quad (4.4.57)$$

for all  $x \in \bar{\Omega}$ ,  $t \geq 0$ . Consequently, the solution  $u(x, t)$  exists for all time.

The region  $\{u \in \mathbb{R}^n : m_\alpha \leq u^\alpha \leq M_\alpha \ (\alpha = 1, \dots, n)\}$  is called an invariant region for the reaction-diffusion system because a solution that starts in it will never leave it. The geometric idea behind this is of course that near the lower boundary value  $m_\alpha$ , the component  $u^\alpha$  of the solution will increase, because of (4.4.55) and the maximum principle Lemma 4.1.2, whereas at the upper value  $M_\alpha$ , it will decrease. In other words, the properties of the reaction terms  $F^\alpha$  force the solution to stay inside the region. Therefore, it has to stay bounded. The precise proofs of these results, however, are too long to be presented here, and we refer to [34, 21] and the references listed there.

We now look at an example where Theorem 4.4.3 can be applied: The FitzHugh-Nagumo equations with diffusion, that is (4.3.201), (4.3.202) for functions  $v(x, t), w(x, t)$  (in place of the notation  $u^1(x, t), u^2(x, t)$ ) that now also depend on a spatial variable

$$v_t = \Delta v + v(a - v)(v - 1) - w \quad (4.4.58)$$

$$w_t = \epsilon \Delta w + bv - cw \quad (4.4.59)$$

for some  $\epsilon \geq 0$ .

We choose  $m_1, M_1, m_2, M_2$  such that (see Figure 4.3.1)

1.  $(m_1, m_2)$  is below the curves  $v(a - v)(v - 1) - w = 0$  and  $bv - cw = 0$  (in particular,  $m_1$  and  $m_2$  are both negative)
2.  $(M_1, m_2)$  is above the curve  $v(a - v)(v - 1) - w = 0$ , but below the curve  $bv - cw = 0$ , and therefore  $bv - cw > 0$  whenever  $m_1 \leq v \leq M_1, w = m_2$
3.  $(m_1, M_2)$  is below  $v(a - v)(v - 1) - w = 0$ , but above the curve  $bv - cw = 0$ ; therefore  $v(a - v)(v - 1) - w > 0$  for  $v = m_1, m_2 \leq w \leq M_2$
4.  $(M_1, M_2)$  is above  $v(a - v)(v - 1) - w = 0$ , and also above  $bv - cw = 0$ ; therefore  $v(a - v)(v - 1) - w < 0$  for  $v = M_1, m_2 \leq w \leq M_2$  as well as  $bv - cw < 0$  whenever  $m_1 \leq v \leq M_1, w = M_2$ .

We observe that in fact we can find arbitrarily large rectangles with these properties. Thus, all assumptions of Theorem 4.4.3 are satisfied for arbitrary bounded initial and boundary values, that is, we can always find an invariant region containing them. We conclude the long-time existence of solutions of (4.4.58), (4.4.59) for any such initial and boundary values.

We now turn to the question of when spatial oscillations die out as time tends to infinity, that is, under which conditions the solution of a reaction-diffusion system tends to a spatially homogeneous state. In order to have access to the simplest situation, in place of the Dirichlet boundary conditions that we have used for our existence results, we now assume homogeneous Neumann boundary conditions

$$\frac{\partial u^\alpha(x, t)}{\partial n} = 0 \text{ for } x \in \partial\Omega, t > 0, \alpha = 1, \dots, n. \quad (4.4.60)$$

For simplicity, we only discuss the case  $F = F(u)$ , that is,  $F$  is independent of  $x$  and  $t$ .

Again, we assume that the solution  $u(x, t)$  stays bounded and consequently exists for all time. We want to compare  $u(x, t)$  with its spatial average  $\bar{u}$  defined by

$$\bar{u}^\alpha(t) := \frac{1}{\|\Omega\|} \int_{\Omega} u^\alpha(x, t) dx \quad (4.4.61)$$

where  $\|\Omega\|$  is the volume of  $\Omega$ .

We also assume a version of (4.4.51):

$$\sup_{x,t} \left\| \frac{dF(u(x, t))}{du} \right\| \leq L. \quad (4.4.62)$$

We let  $\lambda_1 > 0$  be the smallest Neumann eigenvalue of  $\Delta$  on  $\Omega$  (see Theorem 4.1.2).

In order that diffusion can play its role of homogenizing the solution, we need to assume that

$$d := \min_{\alpha=1, \dots, n} d_\alpha > 0 \quad (4.4.63)$$

(this  $d$  should not be confused with the space dimension).

**Theorem 4.4.4.** *Under the assumptions just stated, let  $u(x, t)$  be a bounded solution of (4.4.49) with homogeneous Neumann boundary conditions (4.4.60). If*

$$\delta := d\lambda_1 - L > 0 \quad (4.4.64)$$

*then spatial oscillations of  $u$  decay exponentially on average,*

$$\int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^\alpha(x, t) u_{x^i}^\alpha(x, t) dx \leq c_1 e^{-2\delta t}, \quad (4.4.65)$$

*and  $u$  approaches its spatial average in the  $L^2$ -sense,*

$$\int_{\Omega} |u(x, t) - \bar{u}(t)|^2 dx \leq c_2 e^{-2\delta t}. \quad (4.4.66)$$

*Here,  $c_1, c_2$  are some constants that depend on the various parameters involved.*

*Proof.* The quantity to consider is

$$E(u(\cdot, t)) = \frac{1}{2} \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^\alpha u_{x^i}^\alpha dx.$$

(In more condensed notation,  $E(u(\cdot, t)) = \frac{1}{2} \int_{\Omega} |Du(x, t)|^2 dx$ .)

We compute its temporal evolution:

$$\begin{aligned}
\frac{d}{dt}E(u(\cdot, t)) &= \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{tx^i}^{\alpha} u_{x^i}^{\alpha} dx \\
&= \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^{\alpha} \frac{\partial(d_{\alpha} \Delta u^{\alpha} + F^{\alpha}(u))}{\partial x^i} dx \\
&= - \int_{\Omega} \sum_{\alpha=1}^n d_{\alpha} \Delta u^{\alpha} \Delta u^{\alpha} dx + \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^{\alpha} \sum_{\beta=1}^n \frac{\partial F^{\alpha}}{\partial u^{\beta}} u_{x^i}^{\beta}, \text{ since } \frac{\partial u(x, t)}{\partial \nu} = 0 \text{ for } x \in \partial\Omega \\
&\leq (-d\lambda_1 + L) \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^{\alpha} u_{x^i}^{\alpha} dx = -2\delta E(u(\cdot, t)), \tag{4.4.67}
\end{aligned}$$

using Corollary 4.1.1 and (4.4.64). When we integrate this differential inequality we obtain (4.4.65).

By Corollary 4.1.1, we also have

$$\lambda_1 \int_{\Omega} |u(x, t) - \bar{u}(t)|^2 dx \leq \int_{\Omega} \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^i}^{\alpha}(x, t) u_{x^i}^{\alpha}(x, t) dx, \tag{4.4.68}$$

and therefore, (4.4.65) yields (4.4.66).  $\square$   $\square$

We can also derive a similar result for the temporal variation

**Theorem 4.4.5.** *Under the same assumptions stated, let  $u(x, t)$  again be a bounded solution of (4.4.49) with homogeneous Neumann boundary conditions (4.4.60). If again*

$$\delta = d\lambda_1 - L > 0 \tag{4.4.69}$$

*then temporal oscillations of  $u$  decay exponentially on average,*

$$\int_{\Omega} \sum_{\alpha=1}^n u_t^{\alpha}(x, t) u_t^{\alpha}(x, t) dx \leq c_3 e^{-2\delta t}, \tag{4.4.70}$$

*for some constant  $c_3$ .*

*Proof.* The quantity to consider is

$$E_0(u(\cdot, t)) = \frac{1}{2} \int_{\Omega} \sum_{\alpha=1}^n u_t^{\alpha} u_t^{\alpha} dx.$$

(In more condensed notation,  $E_0(u(\cdot, t)) = \frac{1}{2} \int_{\Omega} |u_t(x, t)|^2 dx$ .)

We compute its temporal evolution:

$$\begin{aligned}
\frac{d}{dt}E_0(u(\cdot, t)) &= \int \sum_{\alpha=1}^n u_t^\alpha u_{tt}^\alpha \\
&= \int \sum_{\alpha=1}^n u_t^\alpha \frac{\partial}{\partial t} (d_\alpha \Delta u^\alpha + F^\alpha(u)) \\
&= -d_\alpha \int \sum_{i=1}^d \sum_{\alpha=1}^n u_{x^{i_t}}^\alpha u_{x^{i_t}}^\alpha + \int \sum_{\alpha=1}^n u_t^\alpha \sum_{\beta=1}^n \frac{\partial F^\alpha}{\partial u^\beta} u_t^\beta \\
&\leq (-\lambda_1 d + L) \int \sum_{\alpha=1}^n u_t^\alpha u_t^\alpha \\
&= 2(-\lambda_1 d + L)E_0(u(\cdot, t)),
\end{aligned}$$

using again Corollary 4.1.1, this time for  $u_t$  which also satisfies a Neumann boundary condition because  $u$  does, and (4.4.69). When we integrate this differential inequality we obtain (4.4.70).  $\square$

When all the diffusion constants  $d_\alpha$  are equal, one can also establish pointwise decay estimates instead of the coarser  $L^2$ -estimates of the preceding theorem.

**Theorem 4.4.6.** *Let  $u(x, t)$  be a bounded solution of*

$$u_t^\alpha(x, t) - \Delta u^\alpha(x, t) = F^\alpha(x, t, u) \quad \text{for } x \in \Omega, t > 0 \quad (4.4.71)$$

*with homogeneous Neumann boundary conditions*

$$\frac{\partial u^\alpha(x, t)}{\partial n} = 0 \quad \text{for } x \in \partial\Omega, t > 0, \alpha = 1, \dots, n. \quad (4.4.72)$$

*If*

$$\delta = \lambda_1 - L > 0 \quad (4.4.73)$$

*then  $u$  approaches its spatial average exponentially,*

$$\sup_{x \in \Omega} |u(x, t) - \bar{u}(t)| \leq c_4 e^{-2\delta t}, \quad (4.4.74)$$

*$c_4$  again being some constant.*

For the *proof*, which needs a stronger analytical tools, namely the regularity theory of parabolic partial differential equations and the Sobolev embedding theorem, to convert integral estimates into pointwise ones, we refer to [21]. Similarly, one may obtain a pointwise decay of  $u_t$ . Thus,  $u$  will tend to a constant as  $t \rightarrow \infty$ , that is, the solution of the reaction-diffusion system will tend towards a homogeneous steady state. Of course, this is not so interesting for pattern formation, and therefore, we now turn to a situation where something else happens.

#### 4.4.4 The Turing mechanism

The Turing mechanism creates instabilities w. r. t. spatial variables for temporally stable states in a system of two coupled reaction-diffusion equations with different diffusion constants.

The system thus is of the form

$$\begin{aligned} A_t &= D_A \Delta A + F(A, B), \\ B_t &= D_B \Delta B + G(A, B). \end{aligned} \quad (4.4.75)$$

*Examples:*

(1) Schnakenberg reaction

$$\begin{aligned} F(A, B) &= k_1 - k_2 A + k_3 A^2 B, \\ G(A, B) &= k_4 - k_3 A^2 B. \end{aligned} \quad (k_{1,2,3,4} > 0)$$

(2) Gierer-Meinhardt system

$$\begin{aligned} F(A, B) &= k_1 - k_2 A + \frac{k_3 A^2}{B}, \\ G(A, B) &= k_4 A^2 - k_5 B. \end{aligned} \quad (k_{1,\dots,5} > 0)$$

(3) Thomas

$$\begin{aligned} F(A, B) &= k_1 - k_2 A - \frac{k_5 AB}{k_6 + k_7 A + k_8 A^2}, \\ G(A, B) &= k_3 - k_4 B - \frac{k_5 AB}{k_6 + k_7 A + k_8 A^2}. \end{aligned} \quad (k_{1,\dots,8} > 0)$$

After rescaling the independent and the dependent variables, the system (4.4.75) becomes

$$\begin{aligned} u_t &= \Delta u + \gamma f(u, v), \\ v_t &= d \Delta v + \gamma g(u, v). \end{aligned} \quad (4.4.76)$$

where the parameter  $\gamma > 0$  is kept for the subsequent analysis.

The preceding examples then become:

(1)

$$\begin{aligned} u_t &= \Delta u + \gamma(a - u + u^2 v), \\ v_t &= d \Delta v + \gamma(b - u^2 v). \end{aligned}$$

(2)

$$\begin{aligned} u_t &= \Delta u + \gamma\left(a - bu + \frac{u^2}{v}\right), \\ v_t &= d\Delta v + \gamma(u^2 - v). \end{aligned}$$

(3)

$$\begin{aligned} u_t &= \Delta u + \gamma\left(a - u - \frac{\rho uv}{1 + u + Ku^2}\right), \\ v_t &= d\Delta v + \gamma\left(\alpha(b - v) - \frac{\rho uv}{1 + u + Ku^2}\right). \end{aligned}$$

A slightly more general version of (2) is

(2')

$$\begin{aligned} u_t &= \Delta u + \gamma\left(a - u + \frac{u^2}{v(1 + ku^2)}\right), \\ v_t &= d\Delta v + \gamma(u^2 - v). \end{aligned}$$

Here  $u, v : \Omega \times \mathbb{R}^+ \rightarrow \mathbb{R}$  for some bounded domain  $\Omega \subset \mathbb{R}^d$ , and we fix the initial values

$$u(x, 0), v(x, 0) \quad \text{for } x \in \Omega,$$

and impose Neumann boundary conditions

$$\frac{\partial u}{\partial n}(x, t) = 0 = \frac{\partial v}{\partial n}(x, t) \quad \text{for all } x \in \partial\Omega, t \geq 0.$$

(We have already seen in 4.4.1 that Neumann boundary conditions are well adapted for comparing the solutions of reaction-diffusion system with the ones of the underlying reaction system. – One can also study periodic boundary conditions, or, more generally, consider  $u, v$  as functions on some compact manifold in place of the domain  $\Omega$ .)

The mechanism starts with a fixed point  $(u_*, v_*)$  of the reaction system:

$$f(u_*, v_*) = 0 = g(u_*, v_*)$$

that is linearly stable. One then investigates its stability under spatially inhomogeneous perturbations. According to the discussion following (4.3.164), the stability as a solution of the reaction system means that all eigenvalues  $\lambda$  of

$$A := \begin{pmatrix} f_u(u_*, v_*) & f_v(u_*, v_*) \\ g_u(u_*, v_*) & g_v(u_*, v_*) \end{pmatrix} \quad (4.4.77)$$

have negative real part,

$$\operatorname{Re}(\lambda) < 0. \quad (4.4.78)$$



These eigenvalues are

$$\lambda_{1,2} = \frac{1}{2}\gamma\left((f_u + g_v) \pm \sqrt{(f_u + g_v)^2 - 4(f_u g_v - f_v g_u)}\right). \quad (4.4.79)$$

where the derivatives of  $f$  and  $g$  are evaluated at  $(u_*, v_*)$ . We have  $\operatorname{Re}(\lambda_1) < 0$  and  $\operatorname{Re}(\lambda_2) < 0$  if

$$f_u + g_v < 0, \quad f_u g_v - f_v g_u > 0. \quad (4.4.80)$$

The linearization of the full reaction-diffusion system about  $(u_*, v_*)$  is

$$w_t = \begin{pmatrix} 1 & 0 \\ 0 & d \end{pmatrix} \Delta w + \gamma A w. \quad (4.4.81)$$

According to Theorem 4.1.2, we let  $u_k$  be an orthonormal basis of eigenfunctions of  $\Delta$  on  $\Omega$  with Neumann boundary conditions,

$$\begin{aligned} \Delta u_k + \lambda_k u_k &= 0 \quad \text{in } \Omega, \\ \frac{\partial u_k}{\partial n} &= 0 \quad \text{on } \partial\Omega. \end{aligned}$$

We then look for solutions of (4.4.81) of the form

$$w_k e^{\lambda t} = \begin{pmatrix} \alpha u_k \\ \beta u_k \end{pmatrix} e^{\lambda t}.$$

Inserting this into (4.4.81) yields

$$\lambda w_k = - \begin{pmatrix} 1 & 0 \\ 0 & d \end{pmatrix} \lambda_k w_k + \gamma A w_k. \quad (4.4.82)$$

For a nontrivial solution of (4.4.82),  $\lambda$  thus has to be an eigenvalue of

$$\left( \gamma A - \begin{pmatrix} 1 & 0 \\ 0 & d \end{pmatrix} \lambda_k \right).$$

The eigenvalue equation is

$$\begin{aligned} \lambda^2 + \lambda(\lambda_k(1+d) - \gamma(f_u + g_v)) \\ + d\lambda_k^2 - \gamma(df_u + g_v)\lambda_k + \gamma^2(f_u g_v - f_v g_u) &= 0. \end{aligned} \quad (4.4.83)$$

We denote the solutions by  $\lambda(k)_{1,2}$ .

(4.4.80) then means that

$$\operatorname{Re} \lambda(0)_{1,2} < 0 \quad (\text{recall } \lambda_0 = 0).$$

We now wish to investigate whether we can have

$$\operatorname{Re} \lambda(k) > 0 \quad (4.4.84)$$

for some higher mode  $\lambda_k$ .

Since by (4.4.80),  $\lambda_k > 0, d > 0$ , clearly

$$\lambda_k(1+d) - \gamma(f_u + g_v) > 0,$$

we need for (4.4.84) that

$$d\lambda_k^2 - \gamma(df_u + g_v)\lambda_k + \gamma^2(f_u g_v - f_v g_u) < 0. \quad (4.4.85)$$

Because of (4.4.80), this can only happen if

$$df_u + g_v > 0.$$

Comparing this with the first equation of (4.4.80), we thus need

$$\begin{aligned} d &> 1, \\ f_u g_v &< 0. \end{aligned}$$

If we assume

$$f_u > 0, \quad g_v < 0, \quad (4.4.86)$$

then we need

$$d > 1. \quad (4.4.87)$$

This is not enough to get (4.4.85) negative. In order to achieve this for some value of  $\lambda_k$ , we first determine that value  $\mu$  of  $\lambda_k$  for which the lhs of (4.4.85) is minimized, i. e.

$$\mu = \frac{\gamma}{2d}(df_u + g_v), \quad (4.4.88)$$

and we then need that the lhs of (4.4.85) becomes negative for  $\lambda_k = \mu$ . This is equivalent to

$$\frac{(df_u + g_v)^2}{4d} > f_u g_v - f_v g_u. \quad (4.4.89)$$

If (4.4.89) holds, then the lhs of (4.4.85) has two values of  $\lambda_k$  where it vanishes, namely

$$\begin{aligned} \mu_{\pm} &= \frac{\gamma}{2d} \left( (df_u + g_v) \pm \sqrt{(df_u + g_v)^2 - 4d(f_u g_v - f_v g_u)} \right) \\ &= \frac{\gamma}{2d} \left( (df_u + g_v) \pm \sqrt{(df_u - g_v)^2 + 4df_v g_u} \right) \end{aligned} \quad (4.4.90)$$

and it becomes negative for

$$\mu_- < \lambda_k < \mu_+. \quad (4.4.91)$$

We conclude

**Lemma 4.4.2.** *Suppose (4.4.89) holds. Then  $(u_*, v_*)$  is spatially unstable w. r. t. the mode  $\lambda_k$ , i. e. there exists a solution of (4.4.82) with*

$$\operatorname{Re} \lambda > 0$$

if  $\lambda_k$  satisfies (4.4.91), where  $\mu_{\pm}$  are given by (4.4.90).

(4.4.89) is satisfied for

$$d > d_c = -\frac{2f_v g_u - f_u g_v}{f_u^2} + \frac{2}{f_u^2} \sqrt{f_v g_u (f_v g_u - f_u g_v)}. \quad (4.4.92)$$

Whether there exists an eigenvalue  $\lambda_k$  of  $\Delta$  satisfying (4.4.91) depends on the geometry of  $\Omega$ . In particular, according to the discussion after Theorem 4.1.2, the eigenvalues scale like  $\|\Omega\|^{-\frac{2}{d}}$ . Thus, if  $\Omega$  is small, all nonzero eigenvalues are large. Therefore, if  $\Omega$  is sufficiently small, all nonzero eigenvalues are larger than  $\mu_+$ .

We can also keep  $\Omega$ , and thus the smallest nonzero eigenvalue  $\lambda_1$ , fixed and choose  $\gamma > 0$  in (4.4.90) so small that

$$\mu_+ < \lambda_1.$$

Then, again, (4.4.91) cannot be solved. From these considerations we see that we need a certain minimal domain size for a given reaction strength, or else a certain minimal reaction strength for a given domain size, for a Turing instability to occur.

If the condition (4.4.91) is satisfied for some eigenvalue  $\lambda_k$ , it is also of geometric significance for which value of  $k$  this happens. Namely, by Courant's nodal domain theorem, the nodal set  $\{u_k = 0\}$  of the eigenfunction  $u_k$  divides  $\Omega$  into at most  $(k+1)$  regions. On any of these regions,  $u_k$  then has a fixed sign, i. e. is either positive or negative on that entire region. Since  $u_k$  is the unstable mode, this controls the number of oscillations of the developing instability.

We summarize our considerations in

**Theorem 4.4.7.** *Suppose that at a solution  $(u_*, v_*)$  of*

$$f(u_*, v_*) = 0 = g(u_*, v_*),$$

we have

$$f_u + g_v < 0, \quad f_u g_v - f_v g_u > 0.$$

Then  $(u_*, v_*)$  is linearly stable for the reaction system

$$\begin{aligned}u_t &= \gamma f(u, v), \\v_t &= \gamma g(u, v).\end{aligned}$$

Suppose that  $d > 1$  satisfies

$$\begin{aligned}df_u + g_v &> 0, \\(df_u + g_v)^2 - 4d(f_u g_v - f_v g_u) &> 0.\end{aligned}$$

Then  $(u_*, v_*)$  as a solution of the reaction-diffusion system

$$\begin{aligned}u_t &= \Delta u + \gamma f(u, v), \\v_t &= d\Delta v + \gamma g(u, v)\end{aligned}$$

is linearly unstable against spatial oscillations with eigenvalue  $\lambda_k$  whenever  $\lambda_k$  satisfies (4.4.91).

Since the eigenvalues  $\lambda_k$  of  $\Delta$  on the bounded domain  $\Omega$  are discrete (recall Theorem 4.1.2), it depends on the geometry of  $\Omega$  whether such an eigenvalue in the range determined by (4.4.91) exists. The number  $k$  controls the frequency of oscillations of the instability about  $(u_*, v_*)$ , and thus determines the shape of the resulting spatial pattern.

Thus, in the situation described in Theorem 4.4.7, the equilibrium state  $(u_*, v_*)$  is unstable, and in the vicinity of it, perturbations grow at a rate  $e^{\text{Re}\lambda}$ , where  $\lambda$  solves (4.4.83).

Typically, one assumes, however, that the dynamics is confined within a bounded region in  $(\mathbb{R}^+)^2$ . This means that appropriate assumptions on  $f$  and  $g$  for  $u = 0$  or  $v = 0$ , or for  $u$  and  $v$  large ensure that solutions starting in the positive quadrant can neither become zero nor unbounded. In 4.4.3 we have discussed the principle that if such a confinement holds for the reaction system, it also holds for the reaction-diffusion system.

Thus, even though  $(u_*, v_*)$  is locally unstable, and therefore small perturbations grow exponentially, this growth has to terminate eventually, and one expects that the corresponding solution of the reaction-diffusion system settles at a spatially inhomogeneous steady state. This is the idea of the Turing mechanism. This has not yet been demonstrated in full rigour and generality. So far, the existence of spatially heterogeneous solutions has only been shown by singular perturbation analysis near the critical parameter  $d_c$  in (4.4.92).

We should also point out that Turing structures show that what we have discussed at the end of 4.4.3, namely that solutions of reaction-diffusion systems become spatially homogeneous as time tends to infinity, is by no means a universal phenomenon, but rather depends on specific assumptions. Clearly, spatially inhomogeneous structures as produced by the Turing mechanism are more interesting for pattern formation than homogeneous ones. Of course, the situation becomes even richer when the asymptotic structures are not only not spatially

homogeneous, but also not steady in time. For example, a Turing instability could get combined with a Hopf bifurcation. See [37] for examples.

With Theorem 4.4.7, we return to example (1) above. We have

$$\begin{aligned} u_* &= a + b, \\ v_* &= \frac{b}{(a + b)^2}, \quad (\text{of course, } a, b > 0) \end{aligned}$$

and at  $(u_*, v_*)$  then

$$\begin{aligned} f_u &= \frac{b - a}{a + b}, \\ f_v &= (a + b)^2, \\ g_u &= -\frac{2b}{a + b}, \\ g_v &= -(a + b)^2, \\ f_u g_v - f_v g_u &= (a + b)^2 > 0. \end{aligned}$$

Since we need that  $f_u$  and  $g_v$  have opposite signs (in order to get  $df_u + g_v > 0$  later on), we require

$$b > a.$$

$f_u + g_v < 0$  then implies

$$0 < b - a < (a + b)^3, \quad (4.4.93)$$

while  $df_u + g_v > 0$  implies

$$d(b - a) > (a + b)^3. \quad (4.4.94)$$

Finally,  $(df_u + g_v)^2 - 4d(f_u g_v - f_v g_u) > 0$  requires

$$(d(b - a) - (a + b)^3)^2 > 4d(a + b)^4. \quad (4.4.95)$$

The Turing mechanism is a beautiful analytical scheme for pattern formation. This, however, does not imply that this really is the general scheme underlying the formation of spatial patterns in biology. In fact, according to recent developments in developmental biology, the combinatorial patterns of gene regulation constitute the basic mechanism for the formation of spatial structures, rather than the Turing mechanism. Nevertheless, in certain cases, Turing's idea [36] may apply. Many examples are discussed in [28].

The present section follows the presentation in [21] rather closely.

## 4.5 Continuity and Fokker-Planck equations

There exists a scheme that represents an alternative to the reaction-diffusion paradigm (4.4.1) for passing from dynamical systems to partial differential equations.

One starts again with a dynamical system

$$u_t(t) = f(t, u). \quad (4.5.96)$$

Here,  $u$  takes its values in  $\mathbb{R}^n$ .

The object of interest then is the density  $h(u, t)$  of  $u$ , that is, for each measurable  $A \subset \mathbb{R}^n$ , the probability that  $u(t)$  is contained in  $A$  is given by

$$\int_A h(y, t) dy. \quad (4.5.97)$$

Of course, this depends only on the initial values  $u(0)$ , and as such, we expect  $h(u, t)$  to evolve as  $\delta(u - u(t))$ , where  $\delta$  is the Dirac functional and  $u(t)$  is the solution of (4.5.96). This formalism becomes more interesting when we consider the simultaneous evolution of a family of initial values, instead of only a single value  $u(0)$ . To make this consistent, we therefore assume that we have an initial density  $h(u, 0)$ , normalized to be a probability, that is,

$$\int_{\mathbb{R}^n} h(y, 0) dy = 1. \quad (4.5.98)$$

We also assume a suitable decay at infinity.

Each value  $u$  evolves according to (4.5.96), and we want to have an equation for the evolution of the density  $h(u, t)$ . This equation is the continuity equation

$$\frac{\partial}{\partial t} h(u, t) = \frac{\partial}{\partial u^i} (-f^i(t, u) h(u, t)). \quad (4.5.99)$$

This equation states that the change of the probability density in time is the negative of the change of the state as a function of its value.

A solution of (4.5.99) then satisfies the normalization (4.5.98) for all time, that is,

$$\int_{\mathbb{R}^n} h(y, t) dy = 1 \text{ for all } t \geq 0. \quad (4.5.100)$$

(4.5.99) is a first order partial differential equation of hyperbolic type. When  $f = 1$ , we obtain a so-called transport equation

$$\frac{\partial}{\partial t} h(u, t) + \frac{\partial}{\partial u} h(u, t). \quad (4.5.101)$$

(4.5.99) represents an important paradigm shift for dynamical systems. So far, we have focussed on individual trajectories, that is, considered the evolution of a single initial value in time. Now, we rather take a family of initial values and investigate how the density of states evolves. The limiting density, if it exists,

would represent a stationary state distribution.

In order to introduce diffusion effects, we now assume that the evolution equation (4.5.96) is subjected to white noise (we recall here the discussion of (4.2.144) above). Formally, one writes

$$u_t(t) = f(t, u) + \eta. \quad (4.5.102)$$

In order to understand what this means, we first put  $f = 0$  and consider

$$u_t(t) = \eta. \quad (4.5.103)$$

This means that  $u$  randomly fluctuates in the sense that

$$u(t) = \int_0^t dw(\tau) + u(0) \quad (4.5.104)$$

where  $w(t)$  is Brownian motion and the integral is a so-called Itô integral. Instead of explaining, however, what that means (see e.g. [22]), we rather state the corresponding equation for the probability density

$$\frac{\partial}{\partial t} h(u, t) = \frac{1}{2} \Delta h(u, t) \quad (4.5.105)$$

where the Laplacian  $\Delta$  operates on the  $u$ -variables, that is,  $\Delta h(u, t) = \sum_{i=1}^n \frac{\partial^2}{(\partial u^i)^2} h(u, t)$ . So, in contrast to the reaction-diffusion paradigm, here, the state variable is not diffusing in physical space, that is, for a variable  $x \in \mathbb{R}^d$ , but rather the state value is randomly fluctuating, leading to a diffusion for its density. – Of course, (4.5.105) is the Fokker-Planck equation already studied above as a continuum limit of Brownian motion.

We can then combine the dynamical system (4.5.96) leading to the continuity equation (4.5.99) and the Fokker-Planck equation (4.5.105), to arrive at

$$\frac{\partial}{\partial t} h(u, t) = \frac{1}{2} \Delta h(u, t) - \frac{\partial}{\partial u} (f(t, u) h(u, t)). \quad (4.5.106)$$

Again, a solution with Neumann boundary conditions satisfies (4.5.100) if it does so for  $t = 0$ .

We should point out the differences between reaction-diffusion equations (4.4.1) or systems on one hand and Fokker-Planck equations on the other hand. (4.5.105) is an equation for a scalar valued quantity  $h$  while the state variable  $u$  can be vector valued, that is  $u \in \mathbb{R}^n$ , and the corresponding  $u$  in a reaction-diffusion system then also takes values in  $\mathbb{R}^n$ . Also, in contrast to reaction-diffusion equations, (4.5.105) does not include physical space. The diffusion takes place in state space, not in physical space. (Of course, within this framework, one can also study the setting where the state space is physical space, that is, where we consider a quantity that moves around in physical space governed by some dynamical law and diffusion.)

For a systematic presentation of modern results on PDE models in mathematical biology, see [31].

Thanks for discussions to Ugur Abdullah, Anirban Banerjee, Andreas Dress, Benoit Perthame, Angela Stevens



# Bibliography

- [1] H.-J. Bandelt, A.Dress, A canonical decomposition theory for metrics on a finite set, *Adv.Math.*92, 47-105, 1992
- [2] E.Bender, E.Canfield, The asymptotic number of labeled graphs with given degree sequences, *J.Comb.Th.A* 24, 296-307, 1978
- [3] B.Bolobás, *Random graphs*, Cambridge Univ.Press, <sup>2</sup>2001
- [4] B.Bolobás, *Modern graph theory*, Springer, 1998
- [5] F.Chung, *Spectral graph theory*, AMS, 1997
- [6] H.Colonius, H.Schulze, Tree structures for proximity data, *Brit.J.Math.Statist.Psychol.*34, 167-180, 1981
- [7] R.Courant, K.Friedrichs, H.Lewy, Über die partiellen Differentialgleichungen der mathematischen Physik, *Math.Ann.*100, 32-74, 1928
- [8] P.Dayan, L.F.Abbott, *Theoretical Neuroscience*, MIT Press, 2001
- [9] A.Dress, Recent results and new problems in phylogenetic combinatorics
- [10] P.Erdős, A.Rényi, On random graphs.I, *Publ.Math.Debrecen* 6, 290-291, 1959
- [11] L.C.Evans, *Partial differential equations*, AMS, 1998
- [12] W.Ewens, *Mathematical population genetics, I. Theoretical introduction*, Springer, <sup>2</sup>2004
- [13] G.Gladwell, E.Davies, J.Leydold, and P.Stadler, Discrete nodal domain theorems, *Lin.Alg.Appl.*336, 2001, 51-60
- [14] C.Godsil, G.Royle, *Algebraic graph theory*, Springer, 2001
- [15] G.Grimmett, D.Stirzacker, *Probability and random processes*, Oxford Univ.Press, <sup>3</sup>2001
- [16] P.Haccou, P.Jagers, V.A.Vatutin, *Branching processes*, Cambridge Univ.Press, 2005

- [17] J.Hein, M.Schierup. C.Wiuf, Gene genealogies, variation and evolution, Oxford Univ.Press, 2005
- [18] W.Hennig, Phylogenetische Systematik, Paul Parey, Berlin and Hamburg, 1982
- [19] J.Hofbauer, K.Sigmund, The theory of evolution and dynamical systems, Cambridge Univ. Press, 1988
- [20] J.Hofbauer, K.Sigmund, Evolutionary games and population dynamics, Cambridge Univ. Press, 1998
- [21] J.Jost, Partial differential equations, Springer, <sup>2</sup>2007
- [22] J.Jost, Dynamical systems, Springer, 2005
- [23] M.Kimmel, D.Axelrod, Branching processes in biology, Springer, 2002
- [24] E.Klipp, R.Herwig, A.Kowald, C.Wierling, H.Lehrach, Systems biology in practice, Wiley-VCH, 2005
- [25] C.Koch, Biophysics of computation, Oxford Univ.Press, New York, 1999
- [26] A.Kolmogoroff, I.Petrovsky, N.Piscounoff, Étude de l' équation de la diffusion avec croissance de la quantité de la matière et son application à un problème biologique, Moscow Univ.Bull.Math.1, 1937, 1-25
- [27] A.Lasota, M.Mackey, Chaos, fractals, and noise, Springer, <sup>2</sup>1994
- [28] J.Murray, Mathematical biology, Springer, 1989
- [29] M.Newman, Random graphs as models of networks, in: S.Bornholdt, H.G.Schuster (eds.), Handbook of graphs and networks, Wiley-VCH, Berlin, 2002
- [30] B.Øksendal, Stochastic differential equations, Springer, <sup>4</sup>1995
- [31] B.Perthame, Transport equations in biology, Birkhäuser, 2007
- [32] P.Phillipson, P.Schuster, An analytical picture of neuron oscillations, Int.J.Bifurc.Chaos 14, 1539-1548, 2004
- [33] C.Semple, M.Steel, Phylogenetics, Oxford Univ.Press, 2003
- [34] J.Smoller, Shock waves and reaction-diffusion equations, Springer, 1983
- [35] K.Strimmer, A. von Haeseler, Quartet puzzling: a quartet maximum likelihood method for reconstructing tree topologies, Mol.BiolEvol.13, 964-996, 1996
- [36] A.Turing
- [37] Walgraef