



## TOWARD A THEORY OF CHAOS

A. SENGUPTA

*Department of Mechanical Engineering,  
 Indian Institute of Technology Kanpur,  
 Kanpur 208016, India  
 osequ@iitk.ac.in*

Received February 23, 2001; Revised September 19, 2002

This paper formulates a new approach to the study of chaos in discrete dynamical systems based on the notions of inverse ill-posed problems, set-valued mappings, generalized and multi-valued inverses, graphical convergence of a net of functions in an extended multifunction space [Sengupta & Ray, 2000], and the topological theory of convergence. Order, chaos and complexity are described as distinct components of this unified mathematical structure that can be viewed as an application of the theory of convergence in topological spaces to increasingly nonlinear mappings, with the boundary between order and complexity in the topology of graphical convergence being the region in  $(\text{Multi}(X))$  that is susceptible to chaos. The paper uses results from the discretized spectral approximation in neutron transport theory [Sengupta, 1988, 1995] and concludes that the numerically exact results obtained by this approximation of the Case singular eigenfunction solution is due to the graphical convergence of the Poisson and conjugate Poisson kernels to the Dirac delta and the principal value multifunctions respectively. In  $(\text{Multi}(X))$ , the continuous spectrum is shown to reduce to a point spectrum, and we introduce a notion of *latent chaotic states* to interpret superposition over generalized eigenfunctions. Along with these latent states, spectral theory of nonlinear operators is used to conclude that nature supports complexity to attain efficiently a multiplicity of states that otherwise would remain unavailable to it.

*Keywords:* Chaos; complexity; ill-posed problems; graphical convergence; topology; multifunctions.

### Prologue

1. *Generally speaking, the analysis of chaos is extremely difficult. While a general definition for chaos applicable to most cases of interest is still lacking, mathematicians agree that for the special case of iteration of transformations there are three common characteristics of chaos:*

1. *Sensitive dependence on initial conditions,*
2. *Mixing,*
3. *Dense periodic points.*

[Peitgen, Jurgens & Saupe, 1992]

2. *The study of chaos is a part of a larger program*

*of study of so-called "strongly" nonlinear system. . . . Linearity means that the rule that determines what a piece of a system is going to do next is not influenced by what it is doing now. More precisely this is intended in a differential or incremental sense: For a linear spring, the increase of its tension is proportional to the increment whereby it is stretched, with the ratio of these increments exactly independent of how much it has already been stretched. Such a spring can be stretched arbitrarily far . . . . Accordingly no real spring is linear. The mathematics of linear objects is particularly felicitous. As it happens, linear objects enjoy an identical, simple geometry. The simplicity of*

*this geometry always allows a relatively easy mental image to capture the essence of a problem, with the technicality, growing with the number of parts, basically a detail. The historical prejudice against nonlinear problems is that no so simple nor universal geometry usually exists.*

Mitchell Feigenbaum's *Foreword* (pp. 1–7)  
in [Peitgen *et al.*, 1992]

3. *The objective of this symposium is to explore the impact of the emerging science of chaos on various disciplines and the broader implications for science and society. The characteristic of chaos is its universality and ubiquity. At this meeting, for example, we have scholars representing mathematics, physics, biology, geophysics and geophysiology, astronomy, medicine, psychology, meteorology, engineering, computer science, economics and social sciences.<sup>1</sup> Having so many disciplines meeting together, of course, involves the risk that we might not always speak the same language, even if all of us have come to talk about "chaos".*

Opening address of Heitor Gurgulino de Souza,  
Rector United Nations University, Tokyo  
[de Souza, 1997]

4. *The predominant approach (of how the different fields of science relate to one other) is reductionist: Questions in physical chemistry can be understood in terms of atomic physics, cell biology in terms of how biomolecules work . . . . We have the best of reasons for taking this reductionist approach: it works. But shortfalls in reductionism are increasingly apparent (and) there is something to be gained from supplementing the predominantly reductionist approach with an integrative agenda. This special section on complex systems is an initial scan (where) we have taken a "complex system" to be one whose properties are not fully explained by an understanding of its component parts. Each Viewpoint author<sup>2</sup> was invited to define "complex" as it applied to his or her discipline.*

[Gallagher & Appenzeller, 1999]

5. *One of the most striking aspects of physics is the simplicity of its laws. Maxwell's equations, Schroedinger's equations, and Hamilton mechanics can each be expressed in a few lines. . . . Everything is simple and neat except, of course, the world. Every place we look outside the physics classroom we see a world of amazing complexity. . . . So why, if the laws are so simple, is the world so complicated? To us complexity means that we have structure with variations. Thus a living organism is complicated because it has many different working parts, each formed by variations in the working out of the same genetic coding. Chaos is also found very frequently. In a chaotic world it is hard to predict which variation will arise in a given place and time. A complex world is interesting because it is highly structured. A chaotic world is interesting because we do not know what is coming next. Our world is both complex and chaotic. Nature can produce complex structures even in simple situations and obey simple laws even in complex situations.*

[Goldenfeld & Kadanoff, 1999]

6. *Where chaos begins, classical science stops. For as long as the world has had physicists inquiring into the laws of nature, it has suffered a special ignorance about disorder in the atmosphere, in the turbulent sea, in the fluctuations in the wildlife populations, in the oscillations of the heart and the brain. But in the 1970s a few scientists began to find a way through disorder. They were mathematicians, physicists, biologists, chemists . . . (and) the insights that emerged led directly into the natural world: the shapes of clouds, the paths of lightning, the microscopic intertwining of blood vessels, the galactic clustering of stars. . . . Chaos breaks across the lines that separate scientific disciplines, (and) has become a shorthand name for a fast growing movement that is reshaping the fabric of the scientific establishment.*

[Gleick, 1987]

7. *order (→) complexity (→) chaos.*

[Waldrop, 1992]

<sup>1</sup>A partial listing of papers is as follows: *Chaos and Politics: Application of Nonlinear Dynamics to Socio-Political issues; Chaos in Society: Reflections on the Impact of Chaos Theory on Sociology; Chaos in Neural Networks; The Impact of Chaos on Mathematics; The Impact of Chaos on Physics; The Impact of Chaos on Economic Theory; The Impact of Chaos on Engineering; The impact of Chaos on Biology; Dynamical Disease: And The Impact of Nonlinear Dynamics and Chaos on Cardiology and Medicine.*

<sup>2</sup>The eight Viewpoint articles are titled: *Simple Lessons from Complexity; Complexity in Chemistry; Complexity in Biological Signaling Systems; Complexity and the Nervous System; Complexity, Pattern, and Evolutionary Trade-Offs in Animal Aggregation; Complexity in Natural Landform Patterns; Complexity and Climate, and Complexity and the Economy.*

8. *Our conclusions based on these examples seem simple: At present chaos is a philosophical term, not a rigorous mathematical term. It may be a subjective notion illustrating the present day limitations of the human intellect or it may describe an intrinsic property of nature such as the “randomness” of the sequence of prime numbers. Moreover, chaos may be undecidable in the sense of Godel in that no matter what definition is given for chaos, there is some example of chaos which cannot be proven to be chaotic from the definition.*

[Brown & Chua, 1996]

9. *My personal feeling is that the definition of a “fractal” should be regarded in the same way as the biologist regards the definition of “life”. There is no hard and fast definition, but just a list of properties characteristic of a living thing . . . . Most living things have most of the characteristics on the list, though there are living objects that are exceptions to each of them. In the same way, it seems best to regard a fractal as a set that has properties such as those listed below, rather than to look for a precise definition which will certainly exclude some interesting cases.*

[Falconer, 1990]

10. *Dynamical systems are often said to exhibit chaos without a precise definition of what this means.*

[Robinson, 1999]

## 1. Introduction

The purpose of this paper is to present a unified, self-contained mathematical structure and physical understanding of the nature of chaos in a discrete dynamical system and to suggest a plausible explanation of *why* natural systems tend to be chaotic. The somewhat extensive quotations with which we begin above, bear testimony to both the increasingly significant — and perhaps all-pervasive — role of nonlinearity in the world today as also our imperfect state of understanding of its manifestations. The list of papers at both the UN Conference [de Souza, 1997] and in *Science* [Gallagher & Appenzeller, 1999] is noteworthy if only to justify the observation of Gleick [1987] that “chaos seems to be everywhere”. Even as everybody appears to be finding chaos and complexity in all likely and unlikely places, and possibly because of it, it is nec-

essary that we have a mathematically clear physical understanding of these notions that are supposedly reshaping our view of nature. This paper is an attempt to contribute to this goal. To make this account essentially self-contained we include here, as far as this is practical, the basics of the background material needed to understand the paper in the form of *Tutorials* and an extended *Appendix*.

The paradigm of chaos of the kneading of the dough is considered to provide an intuitive basis of the mathematics of chaos [Peitgen *et al.*, 1992], and one of our fundamental objectives here is to recount the mathematical framework of this process in terms of the theory of ill-posed problems arising from non-injectivity [Sengupta, 1997], *maximal ill-posedness*, and *graphical convergence* of functions [Sengupta & Ray, 2000]. A natural mathematical formulation of the kneading of the dough in the form of *stretch-cut-and-paste* and *stretch-cut-and-fold* operations is in the ill-posed problem arising from the increasing non-injectivity of the function  $f$  modeling the kneading operation.

---

### ***Begin Tutorial 1: Functions and Multifunctions***

A *relation*, or *correspondence*, between two sets  $X$  and  $Y$ , written  $\mathcal{M}: X \dashrightarrow Y$ , is basically a rule that associates subsets of  $X$  to subsets of  $Y$ ; this is often expressed as  $(A, B) \in \mathcal{M}$  where  $A \subset X$  and  $B \subset Y$  and  $(A, B)$  is an ordered pair of sets. The domain

$$\mathcal{D}(\mathcal{M}) \stackrel{\text{def}}{=} \{A \subset X : (\exists Z \in \mathcal{M})(\pi_X(Z) = A)\}$$

and range

$$\mathcal{R}(\mathcal{M}) \stackrel{\text{def}}{=} \{B \subset Y : (\exists Z \in \mathcal{M})(\pi_Y(Z) = B)\}$$

of  $\mathcal{M}$  are respectively the sets of  $X$  which under  $\mathcal{M}$  corresponds to sets in  $Y$ ; here  $\pi_X$  and  $(\pi_Y)$  are the projections of  $Z$  on  $X$  and  $Y$ , respectively. Equivalently,  $(\mathcal{D}(\mathcal{M}) = \{x \in X : \mathcal{M}(x) \neq \emptyset\})$  and  $(\mathcal{R}(\mathcal{M}) = \bigcup_{x \in \mathcal{D}(\mathcal{M})} \mathcal{M}(x))$ . The *inverse*  $\mathcal{M}^-$  of  $\mathcal{M}$  is the relation

$$\mathcal{M}^- = \{(B, A) : (A, B) \in \mathcal{M}\}$$

so that  $\mathcal{M}^-$  assigns  $A$  to  $B$  iff  $\mathcal{M}$  assigns  $B$  to  $A$ . In general, a relation may assign many elements in its range to a single element from its domain; of

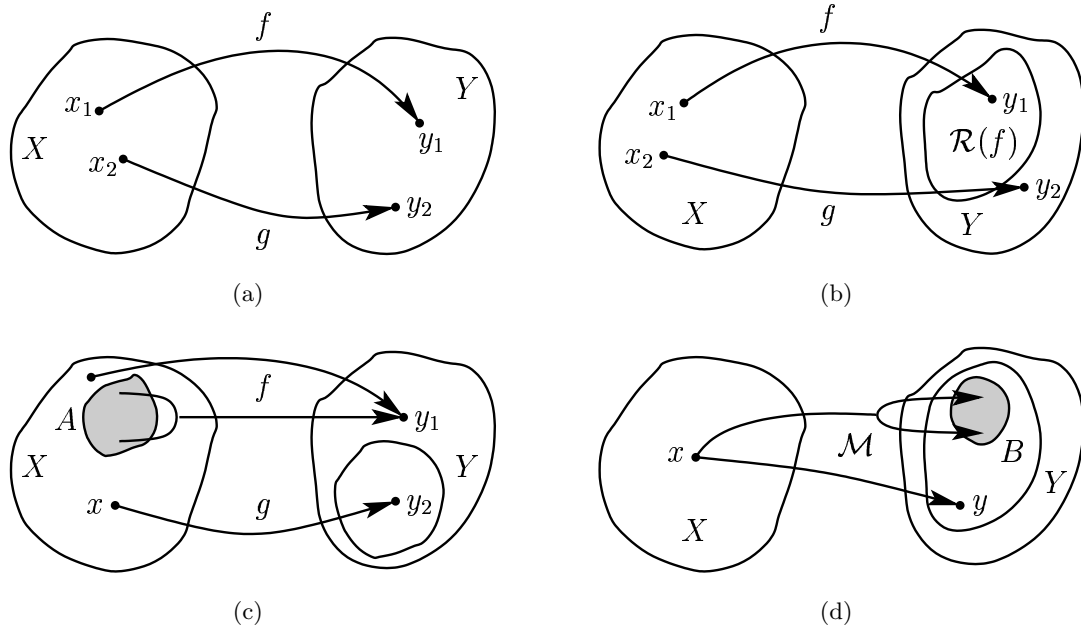


Fig. 1. Functional and non-functional relations between two sets  $X$  and  $Y$ : while  $f$  and  $g$  are functional relations,  $\mathcal{M}$  is not. (a)  $f$  and  $g$  are both injective and surjective (i.e. they are bijective), (b)  $g$  is bijective but  $f$  is only injective and  $f^{-1}(\{y_2\}) := \emptyset$ , (c)  $f$  is not 1:1,  $g$  is not onto, while (d)  $\mathcal{M}$  is not a function but is a *multifunction*.

especial significance are *functional relations*  $f^3$  that can assign only a unique element in  $\mathcal{R}(f)$  to any element in  $\mathcal{D}(f)$ . Figure 1 illustrates the distinction between arbitrary and functional relations  $\mathcal{M}$  and  $f$ . This difference between functions (or maps) and multifunctions is basic to our development and should be fully understood. Functions can again be classified as injections (or 1:1) and surjections (or onto).  $f: X \rightarrow Y$  is said to be *injective* (or *one-to-one*) if  $x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$  for all  $x_1, x_2 \in X$ , while it is *surjective* (or *onto*) if  $Y = f(X)$ .  $f$  is *bijective* if it is both 1:1 and onto.

Associated with a function  $f: X \rightarrow Y$  is its inverse  $f^{-1}: Y \supseteq \mathcal{R}(f) \rightarrow X$  that exists on  $\mathcal{R}(f)$  iff  $f$  is injective. Thus when  $f$  is bijective,  $f^{-1}(y) := \{x \in X: y = f(x)\}$  exists for every  $y \in Y$ ; in fact  $f$  is bijective iff  $f^{-1}(\{y\})$  is a singleton for each  $y \in Y$ . Non-injective functions are not at all rare; if anything, they are very common even for linear maps and it would be perhaps safe to conjecture that they are overwhelmingly predominant in the non-linear world of nature. Thus for example, the simple

linear homogeneous differential equation with constant coefficients of order  $n > 1$  has  $n$  linearly independent solutions so that the operator  $D^n$  of  $D^n(y) = 0$  has a  $n$ -dimensional null space. Inverses of non-injective, and in general non-bijective, functions will be denoted by  $f^-$ . If  $f$  is not injective then

$$A \subset f^-f(A) \stackrel{\text{def}}{=} \text{sat}(A)$$

where  $\text{sat}(A)$  is the *saturation* of  $A \subseteq X$  induced by  $f$ ; if  $f$  is not surjective then

$$ff^-(B) := B \cap f(X) \subseteq B.$$

If  $A = \text{sat}(A)$ , then  $A$  is said to be *saturated*, and  $B \subseteq \mathcal{R}(f)$  whenever  $ff^-(B) = B$ . Thus for non-injective  $f$ ,  $f^-f$  is not an identity on  $X$  just as  $ff^-$  is not  $\mathbf{1}_Y$  if  $f$  is not surjective. However the set of relations

$$ff^-f = f, \quad f^-ff^- = f^- \tag{1}$$

that is always true will be of basic significance in this work. Following are some equivalent statements

<sup>3</sup>We do not distinguish between a relation and its graph although technically they are different objects. Thus although a functional relation, strictly speaking, is the triple  $(X, f, Y)$  written traditionally as  $f: X \rightarrow Y$ , we use it synonymously with the graph  $f$  itself. Parenthetically, the word *functional* in this paper is not necessarily employed for a scalar-valued function, but is used in a wider sense to distinguish between a function and an arbitrary relation (that is a multifunction). Formally, whereas an arbitrary relation from  $X$  to  $Y$  is a subset of  $X \times Y$ , a functional relation must satisfy an additional restriction that requires  $y_1 = y_2$  whenever  $(x, y_1) \in f$  and  $(x, y_2) \in f$ . In this subset notation,  $(x, y) \in f \Leftrightarrow y = f(x)$ .

on the injectivity and surjectivity of functions  $f: X \rightarrow Y$ .

(Injec)  $f$  is 1:1  $\Leftrightarrow$  there is a function  $f_L: Y \rightarrow X$  called the left inverse of  $f$ , such that  $f_L f = \mathbf{1}_X \Leftrightarrow A = f^{-1} f(A)$  for all subsets  $A$  of  $X \Leftrightarrow f(\bigcap A_i) = \bigcap f(A_i)$ .

(Surjec)  $f$  is onto  $\Leftrightarrow$  there is a function  $f_R: Y \rightarrow X$  called the right inverse of  $f$ , such that  $f f_R = \mathbf{1}_Y \Leftrightarrow B = f f^{-1}(B)$  for all subsets  $B$  of  $Y$ .

As we are primarily concerned with non-injectivity of functions, saturated sets generated by equivalence classes of  $f$  will play a significant role in our discussions. A relation  $\mathcal{E}$  on a set  $X$  is said to be an *equivalence relation* if it is<sup>4</sup>

- (ER1) Reflexive:  $(\forall x \in X)(x \mathcal{E} x)$ .
- (ER2) Symmetric:  $(\forall x, y \in X)(x \mathcal{E} y \Rightarrow y \mathcal{E} x)$ .
- (ER3) Transitive:  $(\forall x, y, z \in X)(x \mathcal{E} y \wedge y \mathcal{E} z \Rightarrow x \mathcal{E} z)$ .

Equivalence relations group together unequal elements  $x_1 \neq x_2$  of a set as equivalent according to the requirements of the relation. This is expressed as  $x_1 \sim x_2 \pmod{\mathcal{E}}$  and will be represented here by the shorthand notation  $x_1 \sim_{\mathcal{E}} x_2$ , or even simply as  $x_1 \sim x_2$  if the specification of  $\mathcal{E}$  is not essential. Thus for a non-injective map if  $f(x_1) = f(x_2)$  for  $x_1 \neq x_2$ , then  $x_1$  and  $x_2$  can be considered to be equivalent to each other since they map onto the same point under  $f$ ; thus  $x_1 \sim_f x_2 \Leftrightarrow f(x_1) = f(x_2)$  defines the equivalence relation  $\sim_f$  induced by the map  $f$ . Given an equivalence relation  $\sim$  on a set  $X$  and an element  $x \in X$  the subset

$$[x] \stackrel{\text{def}}{=} \{y \in X : y \sim x\}$$

is called the *equivalence class* of  $x$ ; thus  $x \sim y \Leftrightarrow [x] = [y]$ . In particular, equivalence classes generated by  $f: X \rightarrow Y$ ,  $[x]_f = \{x_\alpha \in X : f(x_\alpha) = f(x)\}$ , will be a cornerstone of our analysis of chaos generated by the iterates of non-injective maps, and the equivalence relation  $\sim_f := \{(x, y) : f(x) = f(y)\}$  generated by  $f$  is uniquely defined by the partition that  $f$  induces on  $X$ . Of course as  $x \sim x$ ,  $x \in [x]$ . It is a simple matter to see that any two equivalence classes are either disjoint or equal so that the equivalence classes generated by an equivalence relation on  $X$  form a disjoint cover of  $X$ . The *quotient*

set of  $X$  under  $\sim$ , denoted by  $X/\sim := \{[x] : x \in X\}$ , has the equivalence classes  $[x]$  as its elements; thus  $[x]$  plays a dual role either as subsets of  $X$  or as elements of  $X/\sim$ . The rule  $x \mapsto [x]$  defines a surjective function  $Q: X \rightarrow X/\sim$  known as the *quotient map*.

**Example 1.1.** Let

$$S^1 = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 = 1\}$$

be the unit circle in  $\mathbb{R}^2$ . Consider  $X = [0, 1]$  as a subspace of  $\mathbb{R}$ , define a map

$$q : X \rightarrow S^1, \quad s \mapsto (\cos 2\pi s, \sin 2\pi s), \quad s \in X,$$

from  $\mathbb{R}$  to  $\mathbb{R}^2$ , and let  $\sim$  be the equivalence relation on  $X$

$$s \sim t \Leftrightarrow (s = t) \vee (s = 0, t = 1) \vee (s = 1, t = 0).$$

If we bend  $X$  around till its ends touch, the resulting circle represents the quotient set  $Y = X/\sim$  whose points are equivalent under  $\sim$  as follows

$$[0] = \{0, 1\} = [1], \quad [s] = \{s\} \text{ for all } s \in (0, 1).$$

Thus  $q$  is bijective for  $s \in (0, 1)$  but two-to-one for the special values  $s = 0$  and  $1$ , so that for  $s, t \in X$ ,

$$s \sim t \Leftrightarrow q(s) = q(t).$$

This yields a bijection  $h: X/\sim \rightarrow S^1$  such that

$$q = h \circ Q$$

defines the quotient map  $Q: X \rightarrow X/\sim$  by  $h([s]) = q(s)$  for all  $s \in [0, 1]$ . The situation is illustrated by the commutative diagram of Fig. 2 that appears as an integral component in a different and more general context in Sec. 2. It is to be noted that commutativity of the diagram implies that if a given equivalence relation  $\sim$  on  $X$  is completely determined by  $q$  that associates the partitioning equivalence classes in  $X$  to unique points in  $S^1$ , then  $\sim$  is identical to the equivalence relation that is induced by  $Q$  on  $X$ . Note that a larger size of the equivalence classes can be obtained by considering  $X = \mathbb{R}_+$  for which  $s \sim t \Leftrightarrow |s - t| \in \mathbb{Z}_+$ . ■

### End Tutorial 1

---

<sup>4</sup>An alternate useful way of expressing these properties for a relation  $\mathcal{R}$  on  $X$  are

- (ER1)  $\mathcal{R}$  is reflexive iff  $\mathbf{1}_X \subseteq \mathcal{R}$
  - (ER2)  $\mathcal{R}$  is symmetric iff  $\mathcal{R} = \mathcal{R}^{-1}$
  - (ER3)  $\mathcal{R}$  is transitive iff  $\mathcal{R} \circ \mathcal{R} \subseteq \mathcal{R}$ ,
- with  $\mathcal{R}$  an equivalence relation only if  $\mathcal{R} \circ \mathcal{R} = \mathcal{R}$ .

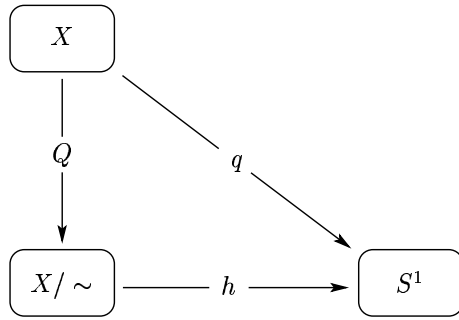


Fig. 2. The quotient map  $Q$ .

One of the central concepts that we consider and employ in this work is the inverse  $f^-$  of a nonlinear, non-injective, function  $f$ ; here the equivalence classes  $[x]_f = f^-f(x)$  of  $x \in X$  are the saturated subsets of  $X$  that partition  $X$ . While a detailed treatment of this question in the form of the nonlinear ill-posed problem and its solution is given in Sec. 2 [Sengupta, 1997], it is sufficient to point out here from Figs. 1(c) and 1(d), that the inverse of a non-injective function is not a function but a multifunction while the inverse of a multifunction is a non-injective function. Hence one has the general result that

$$\begin{aligned}
 & f \text{ is a non-injective function} \\
 & \Leftrightarrow f^- \text{ is a multifunction.} \\
 & f \text{ is a multifunction} \\
 & \Leftrightarrow f^- \text{ is a non-injective function}
 \end{aligned}
 \tag{2}$$

The inverse of a multifunction  $\mathcal{M}: X \multimap Y$  is a generalization of the corresponding notion for a function  $f: X \rightarrow Y$  such that

$$\mathcal{M}^-(y) \stackrel{\text{def}}{=} \{x \in X : y \in \mathcal{M}(x)\}$$

leads to

$$\mathcal{M}^-(B) = \{x \in X : \mathcal{M}(x) \cap B \neq \emptyset\}$$

for any  $B \subseteq Y$ , while a more restricted inverse that we shall not be concerned with is given as  $\mathcal{M}^+(B) = \{x \in X : \mathcal{M}(x) \subseteq B\}$ . Obviously,  $\mathcal{M}^+(B) \subseteq \mathcal{M}^-(B)$ . A multifunction is injective if  $x_1 \neq x_2 \Rightarrow \mathcal{M}(x_1) \cap \mathcal{M}(x_2) = \emptyset$ , and commonly with functions, it is true that  $\mathcal{M}(\bigcup_{\alpha \in \mathbb{D}} A_\alpha) =$

$\bigcup_{\alpha \in \mathbb{D}} \mathcal{M}(A_\alpha)$  and  $\mathcal{M}(\bigcap_{\alpha \in \mathbb{D}} A_\alpha) \subseteq \bigcap_{\alpha \in \mathbb{D}} \mathcal{M}(A_\alpha)$  where  $\mathbb{D}$  is an index set. The following illustrates the difference between the two inverses of  $\mathcal{M}$ . Let  $X$  be a set that is partitioned into two disjoint  $\mathcal{M}$ -invariant subsets  $X_1$  and  $X_2$ . If  $x \in X_1$  (or  $x \in X_2$ ) then  $\mathcal{M}(x)$  represents that part of  $X_1$  (or of  $X_2$ ) that is realized immediately after one application of  $\mathcal{M}$ , while  $\mathcal{M}^-(x)$  denotes the possible precursors of  $x$  in  $X_1$  (or of  $X_2$ ) and  $\mathcal{M}^+(B)$  is that subset of  $X$  whose image lies in  $B$  for any subset  $B \subset X$ .

In this paper the multifunctions that we shall be explicitly concerned with arise as the inverses of non-injective maps.

The second major component of our theory is the graphical convergence of a net of functions to a multifunction. In Tutorial 2 below, we replace for the sake of simplicity and without loss of generality, the net (which is basically a sequence where the index set is not necessarily the positive integers; thus every sequence is a net but the family<sup>5</sup> indexed, for example, by  $\mathbb{Z}$ , the set of all integers, is a net and not a sequence) with a sequence and provide the necessary background and motivation for the concept of graphical convergence.

---

### Begin Tutorial 2: Convergence of Functions

This Tutorial reviews the inadequacy of the usual notions of convergence of functions either to limit functions or to distributions and suggests the motivation and need for introduction of the notion of graphical convergence of functions to multifunctions. Here, we follow closely the exposition of Korevaar [1968], and use the notation  $(f_k)_{k=1}^\infty$  to denote real or complex valued functions on a bounded or unbounded interval  $J$ .

A sequence of piecewise continuous functions  $(f_k)_{k=1}^\infty$  is said to converge to the function  $f$ , notation  $f_k \rightarrow f$ , on a bounded or unbounded interval  $J$ <sup>6</sup>

(1) *Pointwise* if

$$f_k(x) \rightarrow f(x) \quad \text{for all } x \in J,$$

<sup>5</sup>A function  $\chi: \mathbb{D} \rightarrow X$  will be called a *family* in  $X$  indexed by  $\mathbb{D}$  when reference to the domain  $\mathbb{D}$  is of interest, and a *net* when it is required to focus attention on its values in  $X$ .

<sup>6</sup>Observe that it is *not* being claimed that  $f$  belongs to the same class as  $(f_k)$ . This is the single most important cornerstone on which this paper is based: the need to “complete” spaces that are topologically “incomplete”. The classical high-school example of the related problem of having to enlarge, or extend, spaces that are not big enough is the solution space of algebraic equations with real coefficients like  $x^2 + 1 = 0$ .

i.e. Given any arbitrary real number  $\varepsilon > 0$  there exists a  $K \in \mathbb{N}$  that may depend on  $x$ , such that  $|f_k(x) - f(x)| < \varepsilon$  for all  $k \geq K$ .

(2) *Uniformly* if

$$\sup_{x \in J} |f(x) - f_k(x)| \rightarrow 0 \quad \text{as } k \rightarrow \infty,$$

i.e. Given any arbitrary real number  $\varepsilon > 0$  there exists a  $K \in \mathbb{N}$ , such that  $\sup_{x \in J} |f_k(x) - f(x)| < \varepsilon$  for all  $k \geq K$ .

(3) *In the mean of order  $p \geq 1$*  if  $|f(x) - f_k(x)|^p$  is integrable over  $J$  for each  $k$

$$\int_J |f(x) - f_k(x)|^p \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

For  $p = 1$ , this is the simple case of *convergence in the mean*.

(4) *In the mean  $m$ -integrally* if it is possible to select indefinite integrals

$$\begin{aligned} f_k^{(-m)}(x) &= \pi_k(x) + \int_c^x dx_1 \int_c^{x_1} dx_2 \\ &\quad \dots \int_c^{x_{m-1}} dx_m f_k(x_m) \end{aligned}$$

and

$$\begin{aligned} f^{(-m)}(x) &= \pi(x) + \int_c^x dx_1 \int_c^{x_1} dx_2 \\ &\quad \dots \int_c^{x_{m-1}} dx_m f(x_m) \end{aligned}$$

such that for some arbitrary real  $p \geq 1$ ,

$$\int_J |f^{(-m)} - f_k^{(-m)}|^p \rightarrow 0 \quad \text{as } k \rightarrow \infty.$$

where the polynomials  $\pi_k(x)$  and  $\pi(x)$  are of degree  $< m$ , and  $c$  is a constant to be chosen appropriately.

(5) *Relative to test functions  $\varphi$*  if  $f\varphi$  and  $f_k\varphi$  are integrable over  $J$  and

$$\int_J (f_k - f)\varphi \rightarrow 0 \quad \text{for every } \varphi \in \mathcal{C}_0^\infty(J) \text{ as } k \rightarrow \infty,$$

where  $\mathcal{C}_0^\infty(J)$  is the class of infinitely differentiable continuous functions that vanish throughout some neighborhood of each of the end points of  $J$ . For an unbounded  $J$ , a function is said to vanish in some neighborhood of  $+\infty$  if it vanishes on some ray  $(r, \infty)$ .

While pointwise convergence does not imply any other type of convergence, uniform convergence on a bounded interval implies all the other convergences.

It is to be observed that apart from pointwise and uniform convergences, all the other modes listed above represent some sort of an averaged contribution of the entire interval  $J$  and are therefore not of much use when pointwise behavior of the limit  $f$  is necessary. Thus while limits in the mean are not unique, oscillating functions are tamed by  $m$ -integral convergence for adequately large values of  $m$ , and convergence relative to test functions, as we see below, can be essentially reduced to  $m$ -integral convergence. On the contrary, our graphical convergence — which may be considered as a pointwise biconvergence with respect to both the direct and inverse images of  $f$  just as usual pointwise convergence is with respect to its direct image only — allows a sequence (in fact, a net) of functions to converge to an arbitrary relation, unhindered by external influences such as the effects of integrations and test functions. To see how this can indeed matter, consider the following

**Example 1.2.** Let  $f_k(x) = \sin kx$ ,  $k = 1, 2, \dots$  and let  $J$  be any bounded interval of the real line. Then 1-integrally we have

$$f_k^{(-1)}(x) = -\frac{1}{k} \cos kx = -\frac{1}{k} + \int_0^x \sin kx_1 dx_1,$$

which obviously converges to 0 uniformly (and therefore in the mean) as  $k \rightarrow \infty$ . And herein lies the point: even though we cannot conclude about the exact nature of  $\sin kx$  as  $k$  increases indefinitely (except that its oscillations become more and more pronounced), we may very definitely state that  $\lim_{k \rightarrow \infty} (\cos kx)/k = 0$  uniformly. Hence from

$$f_k^{(-1)}(x) \rightarrow 0 = 0 + \int_0^x \lim_{k \rightarrow \infty} \sin kx_1 dx_1$$

it follows that

$$\lim_{k \rightarrow \infty} \sin kx = 0 \tag{3}$$

1-integrally.

Continuing with the same sequence of functions, we now examine its test-functional convergence with respect to  $\varphi \in \mathcal{C}_0^1(-\infty, \infty)$  that vanishes for all  $x \notin (\alpha, \beta)$ . Integrating by parts,

$$\begin{aligned} \int_{-\infty}^\infty f_k \varphi &= \int_\alpha^\beta \varphi(x_1) \sin kx_1 dx_1 \\ &= -\frac{1}{k} [\varphi(x_1) \cos kx_1]_\alpha^\beta \\ &\quad - \frac{1}{k} \int_\alpha^\beta \varphi'(x_1) \cos kx_1 dx_1 \end{aligned}$$

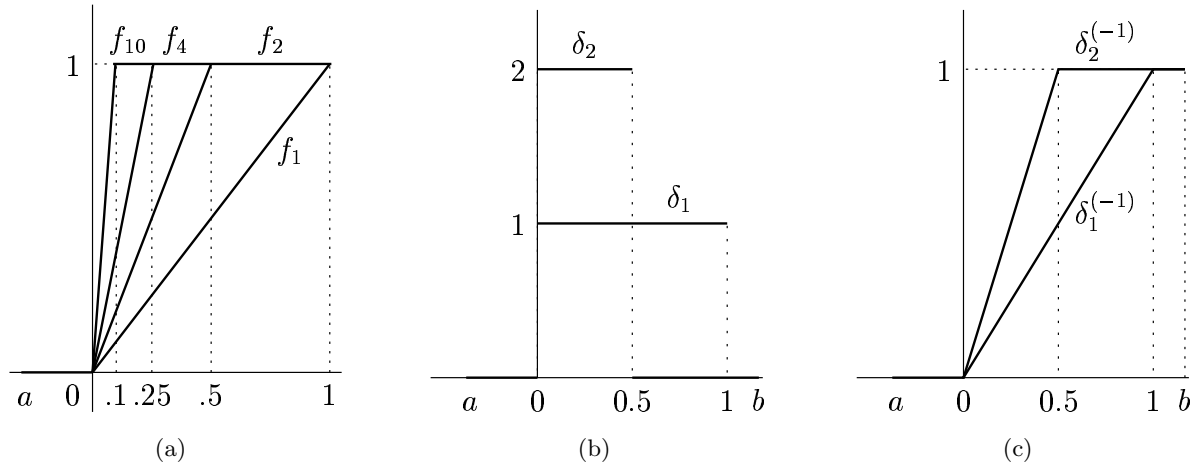


Fig. 3. Incompleteness of function spaces. (a) demonstrates the classic example of non-completeness of the space of real-valued continuous functions leading to the complete spaces  $L_n[a, b]$  whose elements are equivalence classes of functions with  $f \sim g$  iff the Lebesgue integral  $\int_a^b |f - g|^n = 0$ . (b) and (c) illustrate distributional convergence of the functions  $f_k(x)$  of Eq. (5) to the Dirac delta  $\delta(x)$  leading to the complete space of generalized functions. In comparison, note that the space of continuous functions in the uniform metric  $C[a, b]$  is complete which suggests the importance of topologies in determining convergence properties of spaces.

The first integrated term is 0 due to the conditions on  $\varphi$  while the second also vanishes because  $\varphi \in \mathcal{C}_0^1(-\infty, \infty)$ . Hence

$$\int_{-\infty}^{\infty} f_k \varphi \rightarrow 0 = \int_{\alpha}^{\beta} \lim_{k \rightarrow \infty} \varphi(x_1) \sin kx dx_1$$

for all  $\varphi$ , and leading to the conclusion that

$$\lim_{k \rightarrow \infty} \sin kx = 0 \tag{4}$$

test-functionally.

This example illustrates the fact that if  $\text{Supp}(\varphi) = [\alpha, \beta] \subseteq J$ ,<sup>7</sup> integrating by parts sufficiently large number of times so as to wipe out the pathological behavior of  $(f_k)$  gives

$$\begin{aligned} \int_J f_k \varphi &= \int_{\alpha}^{\beta} f_k \varphi \\ &= \int_{\alpha}^{\beta} f_k^{(-1)} \varphi' = \dots = (-1)^m \int_{\alpha}^{\beta} f_k^{(-m)} \varphi^{(m)} \end{aligned}$$

where  $f_k^{(-m)}(x) = \pi_k(x) + \int_c^x dx_1 \int_c^{x_1} dx_2 \dots \int_c^{x_{m-1}} dx_m f_k(x_m)$  is an  $m$ -times arbitrary indefinite integral of  $f_k$ . If now it is true that  $\int_{\alpha}^{\beta} f_k^{(-m)} \rightarrow \int_{\alpha}^{\beta} f^{(-m)}$ , then it must also be true that  $f_k^{(-m)} \varphi^{(m)}$

converges in the mean to  $f^{(-m)} \varphi^{(m)}$  so that

$$\begin{aligned} \int_{\alpha}^{\beta} f_k \varphi &= (-1)^m \int_{\alpha}^{\beta} f_k^{(-m)} \varphi^{(m)} \\ &\rightarrow (-1)^m \int_{\alpha}^{\beta} f^{(-m)} \varphi^{(m)} = \int_{\alpha}^{\beta} f \varphi. \end{aligned}$$

In fact the converse also holds leading to the following Equivalences between  $m$ -convergence in the mean and convergence with respect to test-functions [Korevaar, 1968].

**Type 1 Equivalence.** If  $f$  and  $(f_k)$  are functions on  $J$  that are integrable on every interior subinterval, then the following are equivalent statements.

- (a) For every interior subinterval  $I$  of  $J$  there is an integer  $m_I \geq 0$ , and hence a smallest integer  $m \geq 0$ , such that certain indefinite integrals  $f_k^{(-m)}$  of the functions  $f_k$  converge in the mean on  $I$  to an indefinite integral  $f^{(-m)}$ ; thus  $\int_I |f_k^{(-m)} - f^{(-m)}| \rightarrow 0$ .
- (b)  $\int_J (f_k - f) \varphi \rightarrow 0$  for every  $\varphi \in \mathcal{C}_0^{\infty}(J)$ .

A significant generalization of this Equivalence is obtained by dropping the restriction that the limit object  $f$  be a function. The need for this generalization arises because metric function spaces are

<sup>7</sup>By definition, the support (or supporting interval) of  $\varphi(x) \in \mathcal{C}_0^{\infty}[\alpha, \beta]$  is  $[\alpha, \beta]$  if  $\varphi$  and all its derivatives vanish for  $x \leq \alpha$  and  $x \geq \beta$ .



known not to be complete: Consider the sequence of functions [Fig. 3(a)]

$$f_k(x) = \begin{cases} 0, & \text{if } a \leq x \leq 0 \\ kx, & \text{if } 0 \leq x \leq \frac{1}{k} \\ 1, & \text{if } \frac{1}{k} \leq x \leq b \end{cases} \quad (5)$$

which is not Cauchy in the uniform metric  $\rho(f_j, f_k) = \sup_{a \leq x \leq b} |f_j(x) - f_k(x)|$  but is Cauchy in the mean  $\rho(f_j, f_k) = \int_a^b |f_j(x) - f_k(x)| dx$ , or even pointwise. However in either case,  $(f_k)$  cannot converge in the respective metrics to a *continuous function* and the limit is a discontinuous unit step function

$$\Theta(x) = \begin{cases} 0, & \text{if } a \leq x \leq 0 \\ 1, & \text{if } 0 < x \leq b \end{cases}$$

with graph  $([a, 0], 0) \cup ((0, b], 1)$ , which is also integrable on  $[a, b]$ . Thus even if the limit of the sequence of continuous functions is not continuous, both the limit and the members of the sequence are integrable functions. This Riemann integration is not sufficiently general, however, and this type of integrability needs to be replaced by a much weaker condition resulting in the larger class of the Lebesgue integrable complete space of functions  $L[a, b]$ .<sup>8</sup>

The functions in Fig. 3(b),

$$\delta_k(x) = \begin{cases} k, & \text{if } 0 < x < \frac{1}{k} \\ 0, & \text{if } x \in [a, b] - \left(0, \frac{1}{k}\right), \end{cases}$$

<sup>8</sup>Both Riemann and Lebesgue integrals can be formulated in terms of the so-called *step functions*  $s(x)$ , which are piecewise constant functions with values  $(\sigma_i)_{i=1}^I$  on a finite number of bounded subintervals  $(J_i)_{i=1}^I$  (which may reduce to a point or may not contain one or both of the end points) of a bounded or unbounded interval  $J$ , with integral  $\int_J s(x) dx \stackrel{\text{def}}{=} \sum_{i=1}^I \sigma_i |J_i|$ . While the Riemann integral of a bounded function  $f(x)$  on a bounded interval  $J$  is defined with respect to sequences of step functions  $(s_j)_{j=1}^\infty$  and  $(t_j)_{j=1}^\infty$  satisfying  $s_j(x) \leq f(x) \leq t_j(x)$  on  $J$  with  $\int_J (s_j - t_j) \rightarrow 0$  as  $j \rightarrow \infty$  as  $R \int_J f(x) dx = \lim \int_J s_j(x) dx = \lim \int_J t_j(x) dx$ , the less restrictive Lebesgue integral is defined for arbitrary functions  $f$  over bounded or unbounded intervals  $J$  in terms of Cauchy sequences of step functions  $\int_J |s_i - s_k| \rightarrow 0, i, k \rightarrow \infty$ , converging to  $f(x)$  as

$$s_j(x) \rightarrow f(x) \text{ pointwise almost everywhere on } J,$$

to be

$$\int_J f(x) dx \stackrel{\text{def}}{=} \lim_{j \rightarrow \infty} \int_J s_j(x) dx.$$

That the Lebesgue integral is more general (and therefore is the proper candidate for completion of function spaces) is illustrated by the example of the function defined over  $[0, 1]$  to be 0 on the rationals and 1 on the irrationals for which an application of the definitions verify that while the Riemann integral is undefined, the Lebesgue integral exists and has value 1. The Riemann integral of a bounded function over a bounded interval exists and is equal to its Lebesgue integral. Because it involves a larger family of functions, all integrals in integral convergences are to be understood in the Lebesgue sense.

can be associated with the arbitrary indefinite integrals

$$\Theta_k(x) \stackrel{\text{def}}{=} \delta_k^{(-1)}(x) = \begin{cases} 0, & a \leq x \leq 0 \\ kx, & 0 < x < \frac{1}{k} \\ 1, & \frac{1}{k} \leq x \leq b \end{cases}$$

of Fig. 3(c), which, as noted above, converge in the mean to the unit step function  $\Theta(x)$ ; hence  $\int_{-\infty}^\infty \delta_k \varphi \equiv \int_\alpha^\beta \delta_k \varphi = -\int_\alpha^\beta \delta_k^{(-1)} \varphi' \rightarrow -\int_0^\beta \varphi'(x) dx = \varphi(0)$ . But there can be no *functional relation*  $\delta(x)$  for which  $\int_\alpha^\beta \delta(x) \varphi(x) dx = \varphi(0)$  for all  $\varphi \in C_0^1[\alpha, \beta]$ , so that unlike in the case in Type 1 Equivalence, the limit in the mean  $\Theta(x)$  of the indefinite integrals  $\delta_k^{(-1)}(x)$  cannot be expressed as the indefinite integral  $\delta^{(-1)}(x)$  of some function  $\delta(x)$  on any interval containing the origin. This leads to the second more general type of equivalence.

**Type 2 Equivalence.** If  $(f_k)$  are functions on  $J$  that are integrable on every interior subinterval, then the following are equivalent statements.

- (a) For every interior subinterval  $I$  of  $J$  there is an integer  $m_I \geq 0$ , and hence a smallest integer  $m \geq 0$ , such that certain indefinite integrals  $f_k^{(-m)}$  of the functions  $f_k$  converge in the mean on  $I$  to an integrable function  $\Theta$  which, unlike in Type 1 Equivalence, need not itself be an indefinite integral of some function  $f$ .

(b)  $c_k(\varphi) = \int_J f_k \varphi \rightarrow c(\varphi)$  for every  $\varphi \in \mathcal{C}_0^\infty(J)$ .

Since we are now given that  $\int_I f_k^{(-m)}(x) dx \rightarrow \int_I \Psi(x) dx$ , it must also be true that  $f_k^{(-m)} \varphi^{(m)}$  converges in the mean to  $\Psi \varphi^{(m)}$  whence

$$\begin{aligned} \int_J f_k \varphi &= (-1)^m \int_I f_k^{(-m)} \varphi^{(m)} \\ &\rightarrow (-1)^m \int_I \Psi \varphi^{(m)} \left( \neq (-1)^m \int_I f^{(-m)} \varphi^{(m)} \right). \end{aligned}$$

The natural question that arises at this stage is then: What is the nature of the relation (not function any more)  $\Psi(x)$ ? For this it is now stipulated, despite the non-equality in the equation above, that as in the mean  $m$ -integral convergence of  $(f_k)$  to a function  $f$ ,

$$\Theta(x) := \lim_{k \rightarrow \infty} \delta_k^{(-1)}(x) \stackrel{\text{def}}{=} \int_{-\infty}^x \delta(x') dx' \quad (6)$$

defines the non-functional relation (“generalized function”)  $\delta(x)$  integrally as a solution of the integral equation (6) of the first kind; hence formally<sup>9</sup>

$$\delta(x) = \frac{d\Theta}{dx} \quad (7)$$

## End Tutorial 2

---

The above tells us that the “delta function” is not a function but its indefinite integral is the piecewise continuous function  $\Theta$  obtained as the mean (or pointwise) limit of a sequence of non-differentiable functions with the integral of  $d\Theta_k(x)/dx$  being preserved for all  $k \in \mathbb{Z}_+$ . What then is the delta (and not its integral)? The answer to this question is contained in our multifunctional extension  $\text{Multi}(X, Y)$  of the function space  $\text{Map}(X, Y)$  considered in Sec. 3. Our treatment of ill-posed problems is used to obtain an understanding and interpretation of the numerical results of the discretized spectral approximation in neutron transport theory [Sengupta, 1988, 1995]. The main conclusions are the following: In a one-dimensional discrete system that is governed by the iterates of a nonlinear map, the dynamics is chaotic if and only if the

system evolves to a state of *maximal ill-posedness*. The analysis is based on the non-injectivity, and hence ill-posedness, of the map; this may be viewed as a mathematical formulation of the *stretch-and-fold* and *stretch-cut-and-paste* kneading operations of the dough that are well-established artifacts in the theory of chaos and the concept of maximal ill-posedness helps in obtaining a *physical understanding* of the nature of chaos. We do this through the fundamental concept of the *graphical convergence* of a sequence (generally a net) of functions [Sengupta & Ray, 2000] that is allowed to converge graphically, when the conditions are right, to a set-valued map or multifunction. Since ill-posed problems naturally lead to multifunctional inverses through functional generalized inverses [Sengupta, 1997], it is natural to seek solutions of ill-posed problems in multifunctional space  $\text{Multi}(X, Y)$  rather than in spaces of functions  $\text{Map}(X, Y)$ ; here  $\text{Multi}(X, Y)$  is an extension of  $\text{Map}(X, Y)$  that is generally larger than the smallest dense extension  $\text{Multi}_|(X, Y)$ .

Feedback and iteration are natural processes by which nature evolves itself. Thus almost every process of evolution is a self-correction process by which the system proceeds from the present to the future through a controlled mechanism of input and evaluation of the past. Evolution laws are inherently nonlinear and complex; here *complexity* is to be understood as the natural manifestation of the nonlinear laws that govern the evolution of the system.

This paper presents a mathematical description of complexity based on [Sengupta, 1997] and [Sengupta & Ray, 2000] and is organized as follows. In Sec. 1, we follow [Sengupta, 1997] to give an overview of ill-posed problems and their solution that forms the foundation of our approach. Sections 2 to 4 apply these ideas by defining a chaotic dynamical system as a *maximally ill-posed problem*; by doing this we are able to overcome the limitations of the three Devaney characterizations of chaos [Devaney, 1989] that apply to the specific case of iteration of transformations in a metric space, and the resulting graphical convergence of functions to multifunctions is the basic tool of our approach. Section 5 analyzes graphical convergence in  $\text{Multi}(X)$  for the discretized spectral approximation

---

<sup>9</sup>The observant reader cannot have failed to notice how mathematical ingenuity successfully transferred the “troubles” of  $(\delta_k)_{k=1}^\infty$  to the sufficiently differentiable benevolent receptor  $\varphi$  so as to be able to work backward, via the resultant trouble free  $(\delta_k^{(-m)})_{k=1}^\infty$ , to the final object  $\delta$ . This necessarily hides the true character of  $\delta$  to allow only a view of its integral manifestation on functions. This unfortunately is not general enough in the strongly nonlinear physical situations responsible for chaos, and is the main reason for constructing the multifunctional extension of function spaces that we use.

of neutron transport theory, which suggests a natural link between ill-posed problems and spectral theory of nonlinear operators. This seems to offer an answer to the question of *why* a natural system should increase its complexity, and eventually tend toward chaoticity, by becoming increasingly nonlinear.

## 2. Ill-Posed Problem and Its Solution

This section based on [Sengupta, 1997] presents a formulation and solution of ill-posed problems arising out of the non-injectivity of a function  $f: X \rightarrow Y$  between topological spaces  $X$  and  $Y$ . A workable knowledge of this approach is necessary as our theory of chaos leading to the characterization of chaotic systems as being a *maximally ill-posed* state of a dynamical system is a direct application of these ideas and can be taken to constitute a mathematical representation of the familiar *stretch-cut-and-paste* and *stretch-and-fold* paradigms of chaos. The problem of finding an  $x \in X$  for a given  $y \in Y$  from the functional relation  $f(x) = y$  is an inverse problem that is *ill-posed* (or, the equation  $f(x) = y$  is ill-posed) if any one or more of the following conditions are satisfied.

(IP1)  $f$  is not injective. This *non-uniqueness* problem of the solution for a given  $y$  is the single most significant criterion of ill-posedness used in this work.

(IP2)  $f$  is not surjective. For a  $y \in Y$ , this is the *existence* problem of the given equation.

(IP3) When  $f$  is bijective, the inverse  $f^{-1}$  is not continuous, which means that small changes in  $y$  may lead to large changes in  $x$ .

A problem  $f(x) = y$  for which a solution exists, is unique, and small changes in data  $y$  that lead to only small changes in the solution  $x$  is said to be *well-posed* or *properly posed*. This means that  $f(x) = y$  is well-posed if  $f$  is bijective and the inverse  $f^{-1}: Y \rightarrow X$  is continuous; otherwise the equation is *ill-posed* or *improperly posed*. It is to be noted that the three criteria are not, in general, independent of each other. Thus if  $f$  represents a bijective, bounded linear operator between Banach spaces  $X$  and  $Y$ , then the inverse mapping theorem guarantees that the inverse  $f^{-1}$  is continuous. Hence ill-posedness depends not only on the algebraic structures of  $X, Y, f$  but also on the topologies of  $X$  and  $Y$ .

**Example 2.1.** As a non-trivial example of an inverse problem, consider the heat equation

$$\frac{\partial \theta(x, t)}{\partial t} = c^2 \frac{\partial^2 \theta(x, t)}{\partial x^2}$$

for the temperature distribution  $\theta(x, t)$  of a one-dimensional homogeneous rod of length  $L$  satisfying the initial condition  $\theta(x, 0) = \theta_0(x), 0 \leq x \leq L$ , and boundary conditions  $\theta(0, t) = 0 = \theta(L, t), 0 \leq t \leq T$ , having the Fourier sine-series solution

$$\theta(x, t) = \sum_{n=1}^{\infty} A_n \sin\left(\frac{n\pi}{L}x\right) e^{-\lambda_n^2 t} \tag{8}$$

where  $\lambda_n = (c\pi/a)n$  and

$$A_n = \frac{2}{L} \int_0^a \theta_0(x') \sin\left(\frac{n\pi}{L}x'\right) dx'$$

are the Fourier expansion coefficients. While the direct problem evaluates  $\theta(x, t)$  from the differential equation and initial temperature distribution  $\theta_0(x)$ , the inverse problem calculates  $\theta_0(x)$  from the integral equation

$$\theta_T(x) = \frac{2}{L} \int_0^a k(x, x') \theta_0(x') dx', \quad 0 \leq x \leq L,$$

when this final temperature  $\theta_T$  is known, and

$$k(x, x') = \sum_{n=1}^{\infty} \sin\left(\frac{n\pi}{L}x\right) \sin\left(\frac{n\pi}{L}x'\right) e^{-\lambda_n^2 T}$$

is the kernel of the integral equation. In terms of the final temperature the distribution becomes

$$\theta_T(x) = \sum_{n=1}^{\infty} B_n \sin\left(\frac{n\pi}{L}x\right) e^{-\lambda_n^2 (t-T)} \tag{9}$$

with Fourier coefficients

$$B_n = \frac{2}{L} \int_0^a \theta_T(x') \sin\left(\frac{n\pi}{L}x'\right) dx'.$$

In  $L^2[0, a]$ , Eqs. (8) and (9) at  $t = T$  and  $t = 0$  yield respectively

$$\|\theta_T(x)\|^2 = \frac{L}{2} \sum_{n=1}^{\infty} A_n^2 e^{-2\lambda_n^2 T} \leq e^{-2\lambda_1^2 T} \|\theta_0\|^2 \tag{10}$$

$$\|\theta_0\|^2 = \frac{L}{2} \sum_{n=1}^{\infty} B_n^2 e^{2\lambda_n^2 T}. \tag{11}$$

The last two equations differ from each other in the significant respect that whereas Eq. (10) shows

that the direct problem is well-posed according to (IP3), Eq. (11) means that in the absence of similar bounds the inverse problem is ill-posed.<sup>10</sup>

**Example 2.2.** Consider the Volterra integral equation of the first kind

$$y(x) = \int_a^x r(x')dx' = Kr$$

where  $y, r \in C[a, b]$  and  $K: C[0, 1] \rightarrow C[0, 1]$  is the corresponding integral operator. Since the differential operator  $D = d/dx$  under the sup-norm  $\|r\| = \sup_{0 \leq x \leq 1} |r(x)|$  is unbounded, the inverse problem  $r = Dy$  for a differentiable function  $y$  on  $[a, b]$  is ill-posed, see Example 6.1. However,  $y = Kr$  becomes well-posed if  $y$  is considered to be in  $C^1[0, 1]$  with norm  $\|y\| = \sup_{0 \leq x \leq 1} |Dy|$ . This illustrates the importance of the topologies of  $X$  and  $Y$  in determining the ill-posed nature of the problem when this is due to (IP3).

Ill-posed problems in nonlinear mathematics of type (IP1) arising from the non-injectivity of  $f$  can be considered to be a generalization of non-uniqueness of solutions of linear equations as, for example, in eigenvalue problems or in the solution of a system of linear algebraic equations with a larger number of unknowns than the number of equations. In both cases, for a given  $y \in Y$ , the solution set of the equation  $f(x) = y$  is given by

$$f^-(y) = [x]_f = \{x' \in X : f(x') = f(x) = y\}.$$

A significant point of difference between linear and nonlinear problems is that unlike the special importance of 0 in linear mathematics, there are no preferred elements in nonlinear problems; this leads to a shift of emphasis from the null space of linear problems to equivalence classes for nonlinear equations. To motivate the role of equivalence classes, let us consider the null spaces in the following linear problems.

(a) Let  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  be defined by  $f(x, y) = x + y$ ,  $(x, y) \in \mathbb{R}^2$ . The null space of  $f$  is generated by the equation  $y = -x$  on the  $x$ - $y$  plane, and the graph of  $f$  is the plane passing through the lines  $\rho = x$  and  $\rho = y$ . For each  $\rho \in \mathbb{R}$  the equivalence classes  $f^-(\rho) = \{(x, y) \in \mathbb{R}^2: x + y = \rho\}$  are lines on the graph parallel to the null set.

(b) For a linear operator  $A: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $m < n$ , satisfying (1) and (2), the problem  $Ax = y$  reduces  $A$  to echelon form with rank  $r$  less than  $\min\{m, n\}$ , when the given equations are consistent. The solution however, produces a generalized inverse leading to a set-valued inverse  $A^-$  of  $A$  for which the inverse images of  $y \in \mathcal{R}(A)$  are multivalued because of the non-trivial null space of  $A$  introduced by assumption (1). Specifically, a null-space of dimension  $n - r$  is generated by the free variables  $\{x_j\}_{j=r+1}^n$  which are arbitrary: this is illposedness of type (1). In addition,  $m - r$  rows of the row reduced echelon form of  $A$  have all 0 entries that introduce restrictions on  $m - r$  coordinates  $\{y_i\}_{i=r+1}^m$  of  $y$  which are now related to  $\{y_i\}_{i=1}^r$ : this illustrates ill-posedness of type (2). Inverse ill-posed problems therefore generate multivalued solutions through a generalized inverse of the mapping.

(c) The eigenvalue problem

$$\left(\frac{d^2}{dx^2} + \lambda^2\right)y = 0 \quad y(0) = 0 = y(1)$$

has the following equivalence class of 0

$$[0]_{D^2} = \{\sin(\pi mx)\}_{m=0}^\infty, \quad D^2 = \left(\frac{d^2}{dx^2} + \lambda^2\right),$$

as its eigenfunctions corresponding to the eigenvalues  $\lambda_m = \pi m$ .

Ill-posed problems are primarily of interest to us explicitly as non-injective maps  $f$ , that is under the condition of (IP1). The two other conditions (IP2) and (IP3) are not as significant and play only an implicit role in the theory. In its application to iterative systems, the degree of non-injectivity of  $f$  defined as the number of its injective branches, increases with iteration of the map. A necessary (but not sufficient) condition for chaos to occur is the increasing non-injectivity of  $f$  that is expressed descriptively in the chaos literature as *stretch-and-fold* or *stretch-cut-and-paste* operations. This increasing non-injectivity that we discuss in the following sections, is what causes a dynamical system to tend toward chaoticity. Ill-posedness arising from non-surjectivity of (injective)  $f$  in the form of *regularization* [Tikhonov & Arsenin, 1977] has received wide attention in the literature of ill-posed problems; this however is not of much significance in our work.

<sup>10</sup>Recall that for a linear operator continuity and boundedness are equivalent concepts.

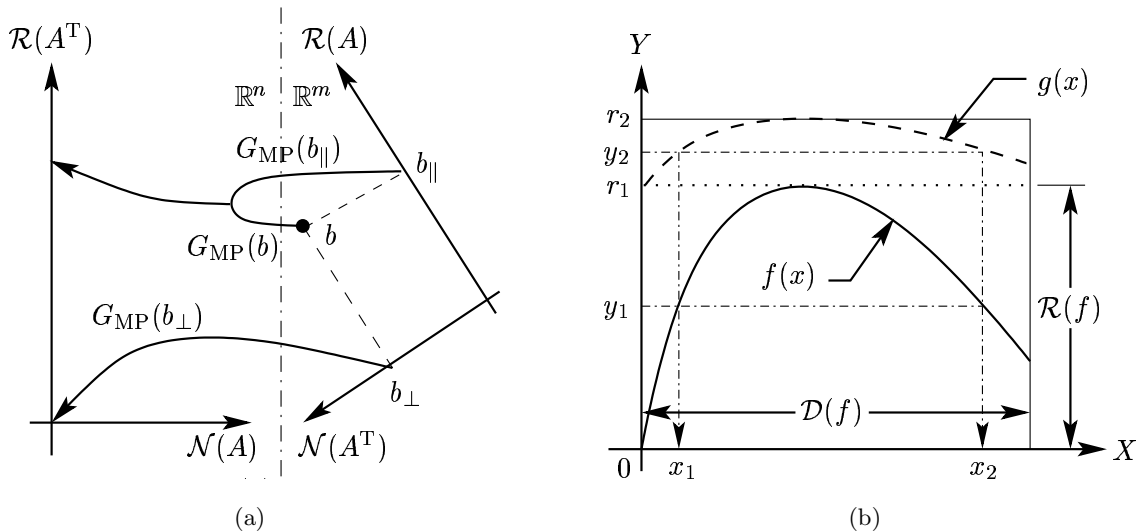


Fig. 4. (a) Moore–Penrose generalized inverse. The decomposition of  $X$  and  $Y$  into the four fundamental subspaces of  $A$  comprising the null space  $\mathcal{N}(A)$ , the column (or range) space  $\mathcal{R}(A)$ , the row space  $\mathcal{R}(A^T)$  and  $\mathcal{N}(A^T)$ , the complement of  $\mathcal{R}(A)$  in  $Y$ , is a basic result in the theory of linear equations. The Moore–Penrose inverse takes advantage of the geometric orthogonality of the row space  $\mathcal{R}(A^T)$  and  $\mathcal{N}(A)$  in  $\mathbb{R}^n$  and that of the column space and  $\mathcal{N}(A^T)$  in  $\mathbb{R}^m$ . (b) When  $X$  and  $Y$  are not inner-product spaces, a *non-injective inverse* can be defined by extending  $f$  to  $Y - \mathcal{R}(f)$  suitably as shown by the dashed curve, where  $g(x) := r_1 + ((r_2 - r_1)/r_1)f(x)$  for all  $x \in \mathcal{D}(f)$  was taken to be a good definition of an extension that replicates  $f$  in  $Y - \mathcal{R}(f)$ ; here  $x_1 \sim x_2$  under both  $f$  and  $g$ , and  $y_1 \sim y_2$  under  $\{f, g\}$  just as  $b$  is equivalent to  $b_{\parallel}$  in the Moore–Penrose case. Note that both  $\{f, g\}$  and  $\{f^-, g^-\}$  are both multifunctions on  $X$  and  $Y$ , respectively. Our inverse  $G$ , introduced later in this section, is however injective with  $G(Y - \mathcal{R}(f)) := 0$ .

### Begin Tutorial 3: Generalized Inverse

In this Tutorial, we take a quick look at the equation  $a(x) = y$ , where  $a: X \rightarrow Y$  is a linear map that need not be either one-one or onto. Specifically, we will take  $X$  and  $Y$  to be the Euclidean spaces  $\mathbb{R}^n$  and  $\mathbb{R}^m$  so that  $a$  has a matrix representation  $A \in \mathbb{R}^{m \times n}$  where  $\mathbb{R}^{m \times n}$  is the collection of  $m \times n$  matrices with real entries. The inverse  $A^{-1}$  exists and is unique iff  $m = n$  and  $\text{rank}(A) = n$ ; this is the situation depicted in Fig. 1(a). If  $A$  is neither one-one or onto, then we need to consider the multifunction  $A^{-}$ , a functional choice of which is known as the *generalized inverse*  $G$  of  $A$ . A good introductory text for generalized inverses is [Campbell & Mayer, 1979]. Figure 4(a) introduces the following definition of the *Moore–Penrose* generalized inverse  $G_{\text{MP}}$ .

**Definition 2.1** (Moore–Penrose Inverse). If  $a: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a linear transformation with matrix representation  $A \in \mathbb{R}^{m \times n}$  then the Moore–Penrose inverse  $G_{\text{MP}} \in \mathbb{R}^{n \times m}$  of  $A$  (we will use the same notation  $G_{\text{MP}}: \mathbb{R}^m \rightarrow \mathbb{R}^n$  for the inverse of the

map  $a$ ) is the noninjective map defined in terms of the row and column spaces of  $A$ ,  $\text{row}(A) = \mathcal{R}(A^T)$ ,  $\text{col}(A) = \mathcal{R}(A)$ , as

$$G_{\text{MP}}(y) \stackrel{\text{def}}{=} \begin{cases} (a|_{\text{row}(A)})^{-1}(y), & \text{if } y \in \text{col}(A) \\ 0, & \text{if } y \in \mathcal{N}(A^T). \end{cases} \tag{12}$$

Note that the restriction  $a|_{\text{row}(A)}$  of  $a$  to  $\mathcal{R}(A^T)$  is bijective so that the inverse  $(a|_{\text{row}(A)})^{-1}$  is well-defined. The role of the transpose matrix appears naturally, and the  $G_{\text{MP}}$  of Eq. (12) is the unique matrix that satisfies the conditions

$$\begin{aligned} AG_{\text{MP}}A &= A, & G_{\text{MP}}AG_{\text{MP}} &= G_{\text{MP}}, \\ (G_{\text{MP}}A)^T &= G_{\text{MP}}A, & (AG_{\text{MP}})^T &= AG_{\text{MP}} \end{aligned} \tag{13}$$

that follow immediately from the definition (12); hence  $G_{\text{MP}}A$  and  $AG_{\text{MP}}$  are orthogonal projections<sup>11</sup> onto the subspaces  $\mathcal{R}(A^T) = \mathcal{R}(G_{\text{MP}})$  and  $\mathcal{R}(A)$ , respectively. Recall that the range space  $\mathcal{R}(A^T)$  of  $A^T$  is the same as the *row space*  $\text{row}(A)$  of  $A$ , and  $\mathcal{R}(A)$  is also known as the *column space* of  $A$ ,  $\text{col}(A)$ .

<sup>11</sup>A real matrix  $A$  is an orthogonal projector iff  $A^2 = A$  and  $A = A^T$ .

**Example 2.3.** For  $a: \mathbb{R}^5 \rightarrow \mathbb{R}^4$ , let

$$A = \begin{pmatrix} 1 & -3 & 2 & 1 & 2 \\ 3 & -9 & 10 & 2 & 9 \\ 2 & -6 & 4 & 2 & 4 \\ 2 & -6 & 8 & 1 & 7 \end{pmatrix}$$

By reducing the augmented matrix  $(A|y)$  to the row-reduced echelon form, it can be verified that the null and range spaces of  $A$  are three- and two-dimensional, respectively. A basis for the null space of  $A^T$  and of the row and column space of  $A$  obtained from the echelon form are respectively

$$\begin{pmatrix} -2 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ 0 \\ 1 \end{pmatrix}; \text{ and } \begin{pmatrix} 1 \\ -3 \\ 0 \\ \frac{3}{2} \\ \frac{1}{2} \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \\ -\frac{1}{4} \\ \frac{3}{4} \end{pmatrix}; \begin{pmatrix} 1 \\ 0 \\ 2 \\ -1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}.$$

According to its definition Eq. (12), the Moore–Penrose inverse maps the middle two of the above set to  $(0, 0, 0, 0, 0)^T$ , and the  $A$ -image of the first two (which are respectively  $(19, 70, 38, 51)^T$  and  $(70, 275, 140, 205)^T$  lying, as they must, in the span of the last two), to the span of  $(1, -3, 2, 1, 2)^T$  and  $(3, -9, 10, 2, 9)^T$  because  $a$  restricted to this subspace of  $\mathbb{R}^5$  is bijective. Hence

$$G_{\text{MP}} \begin{pmatrix} A & \begin{pmatrix} 1 \\ -3 \\ 0 \\ \frac{3}{2} \\ \frac{1}{2} \end{pmatrix} & A & \begin{pmatrix} 0 \\ 0 \\ 1 \\ -\frac{1}{4} \\ \frac{3}{4} \end{pmatrix} & \begin{matrix} -2 & 1 \\ 0 & -1 \\ 1 & 0 \\ 0 & 1 \end{matrix} \end{pmatrix} \\ = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -3 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ \frac{3}{2} & -\frac{1}{4} & 0 & 0 \\ \frac{1}{2} & \frac{3}{4} & 0 & 0 \end{pmatrix}.$$

The second matrix on the left is invertible as its

rank is 4. This gives

$$G_{\text{MP}} = \begin{pmatrix} \frac{9}{275} & -\frac{1}{275} & \frac{18}{275} & -\frac{2}{55} \\ \frac{27}{275} & \frac{3}{275} & \frac{54}{275} & \frac{6}{55} \\ -\frac{10}{143} & \frac{6}{143} & -\frac{20}{143} & \frac{16}{143} \\ \frac{238}{3575} & -\frac{57}{3575} & \frac{476}{3575} & -\frac{59}{715} \\ \frac{129}{3575} & \frac{106}{3575} & \frac{258}{3575} & \frac{47}{715} \end{pmatrix} \tag{14}$$

as the Moore–Penrose inverse of  $A$  that readily verifies all the four conditions of Eqs. (13). The basic point here is that, as in the case of a bijective map,  $G_{\text{MP}}A$  and  $AG_{\text{MP}}$  are identities on the row and column spaces of  $A$  that define its rank. For later use — when we return to this example for a simpler inverse  $G$  — given below are the orthonormal bases of the four fundamental subspaces with respect to which  $G_{\text{MP}}$  is a representation of the generalized inverse of  $A$ ; these calculations were done by MATLAB. The basis for

- (a) the column space of  $A$  consists of the first two columns of the eigenvectors of  $AA^T$ :

$$\begin{pmatrix} -\frac{1633}{2585}, -\frac{363}{892}, \frac{3317}{6387}, \frac{363}{892} \end{pmatrix}^T \\ \begin{pmatrix} -\frac{929}{1435}, \frac{709}{1319}, \frac{346}{6299}, -\frac{709}{1319} \end{pmatrix}^T$$

- (b) the null space of  $A^T$  consists of the last two columns of the eigenvectors of  $AA^T$ :

$$\begin{pmatrix} -\frac{3185}{8306}, \frac{293}{2493}, -\frac{3185}{4153}, \frac{1777}{3547} \end{pmatrix}^T \\ \begin{pmatrix} \frac{323}{1732}, \frac{533}{731}, \frac{323}{866}, \frac{1037}{1911} \end{pmatrix}^T$$

- (c) the row space of  $A$  consists of the first two columns of the eigenvectors of  $A^T A$ :

$$\begin{pmatrix} \frac{421}{13823}, \frac{44}{14895}, -\frac{569}{918}, -\frac{659}{2526}, \frac{1036}{1401} \end{pmatrix} \\ \begin{pmatrix} \frac{661}{690}, \frac{412}{1775}, \frac{59}{2960}, -\frac{1523}{10221}, -\frac{303}{3974} \end{pmatrix}$$

- (d) the null space of  $A$  consists of the last three columns of  $A^T A$ :

$$\begin{pmatrix} -\frac{571}{15469}, & -\frac{369}{776}, & \frac{149}{25344}, & -\frac{291}{350}, & -\frac{389}{1365} \end{pmatrix}$$

$$\begin{pmatrix} -\frac{281}{1313}, & \frac{956}{1489}, & \frac{875}{1706}, & -\frac{1279}{2847}, & \frac{409}{1473} \end{pmatrix}$$

$$\begin{pmatrix} \frac{292}{1579}, & -\frac{876}{1579}, & \frac{203}{342}, & \frac{621}{4814}, & \frac{1157}{2152} \end{pmatrix}$$

The matrices  $Q_1$  and  $Q_2$  with these eigenvectors  $(x_i)$  satisfying  $\|x_i\| = 1$  and  $(x_i, x_j) = 0$  for  $i \neq j$  as their columns are *orthogonal matrices* with the simple inverse criterion  $Q^{-1} = Q^T$ .

**End Tutorial 3**

---

The basic issue in the solution of the inverse ill-posed problem is its reduction to an well-posed one when restricted to suitable subspaces of the domain and range of  $A$ . Considerations of geometry leading to their decomposition into orthogonal subspaces is only an additional feature that is not central to the problem: recall from Eq. (1) that any function  $f$  must necessarily satisfy the more general set-theoretic relations  $ff^{-}f = f$  and  $f^{-}ff^{-} = f^{-}$  of Eq. (13) for the multiinverse  $f^{-}$  of  $f: X \rightarrow Y$ . The second distinguishing feature of the MP-inverse is that it is defined, by a suitable extension, on all of  $Y$  and not just on  $f(X)$  which is perhaps more natural. The availability of orthogonality in inner-product spaces allows this extension to be made in an almost normal fashion. As we shall see below the additional geometric restriction of Eq. (13) is not essential to the solution process, and in fact, only results in a less canonical form of the inverse.

**Begin Tutorial 4: Topological Spaces**

This Tutorial is meant to familiarize the reader with the basic principles of a topological space. A topological space  $(X, \mathcal{U})$  is a set  $X$  with a class<sup>12</sup>  $\mathcal{U}$  of distinguished subsets, called *open sets of  $X$* , that satisfy

- (T1) The empty set  $\emptyset$  and the whole  $X$  belong to  $\mathcal{U}$
- (T2) Finite intersections of members of  $\mathcal{U}$  belong to  $\mathcal{U}$

- (T3) Arbitrary unions of members of  $\mathcal{U}$  belong to  $\mathcal{U}$ .

**Example 2.4**

- (1) The smallest topology possible on a set  $X$  is its *indiscrete topology* when the only open sets are  $\emptyset$  and  $X$ ; the largest is the *discrete topology* where every subset of  $X$  is open (and hence also closed).
- (2) In a metric space  $(X, d)$ , let  $B_\varepsilon(x, d) = \{y \in X: d(x, y) < \varepsilon\}$  be an open ball at  $x$ . Any subset  $U$  of  $X$  such that for each  $x \in U$  there is a  $d$ -ball  $B_\varepsilon(x, d) \subseteq U$  in  $U$ , is said to be an open set of  $(X, d)$ . The collection of all these sets is the topology induced by  $d$ . The topological space  $(X, \mathcal{U})$  is then said to be *associated with (induced by)  $(X, d)$* .
- (3) If  $\sim$  is an equivalence relation on a set  $X$ , the set of all saturated sets  $[x]_\sim = \{y \in X: y \sim x\}$  is a topology on  $X$ ; this topology is called the *topology of saturated sets*.

We argue in Sec. 4.2 that this constitutes the defining topology of a chaotic system.

- (4) For any subset  $A$  of the set  $X$ , the *A-inclusion topology on  $X$*  consists of  $\emptyset$  and every superset of  $A$ , while the *A-exclusion topology on  $X$*  consists of all subsets of  $X - A$ . Thus  $A$  is open in the inclusion topology and closed in the exclusion, and in general every open set of one is closed in the other.

The special cases of the *a-inclusion* and *a-exclusion* topologies for  $A = \{a\}$  are defined in a similar fashion.

- (5) The *cofinite* and *cocountable topologies* in which the open sets of an infinite (resp. uncountable) set  $X$  are respectively the complements of finite and countable subsets, are examples of topologies with some unusual properties that are covered in Appendix A.1. If  $X$  is itself finite (respectively, countable), then its cofinite (respectively, cocountable) topology is the discrete topology consisting of all its subsets. It is therefore useful to adopt the convention, unless stated to the contrary, that cofinite and cocountable spaces are respectively infinite and uncountable.

In the space  $(X, \mathcal{U})$ , a *neighborhood of a point  $x \in X$*  is a nonempty subset  $N$  of  $X$  that contains an open set  $U$  containing  $x$ ; thus  $N \subseteq X$  is a

---

<sup>12</sup>In this sense, a *class* is a set of sets.

neighborhood of  $x$  iff

$$x \in U \subseteq N \tag{15}$$

for some  $U \in \mathcal{U}$ . The largest open set that can be used here is  $\text{Int}(N)$  (where, by definition,  $\text{Int}(A)$  is the largest open set that is contained in  $A$ ) so that the above neighborhood criterion for a subset  $N$  of  $X$  can be expressed in the equivalent form

$$N \subseteq X \text{ is a } \mathcal{U} - \text{neighborhood of } x \text{ iff } x \in \text{Int}_{\mathcal{U}}(N) \tag{16}$$

implying that a subset of  $(X, \mathcal{U})$  is a neighborhood of all its interior points, so that  $N \in \mathcal{N}_x \Rightarrow N \in \mathcal{N}_y$  for all  $y \in \text{Int}(N)$ . The collection of all neighborhoods of  $x$

$$\mathcal{N}_x \stackrel{\text{def}}{=} \{N \subseteq X : x \in U \subseteq N \text{ for some } U \in \mathcal{U}\} \tag{17}$$

is the *neighborhood system* at  $x$ , and the subcollection  $U$  of the topology used in this equation constitutes a *neighborhood (local) base* or *basic neighborhood system, at  $x$* , see Definition A.1.1 of Appendix A.1. The properties

- (N1)  $x$  belongs to every member  $N$  of  $\mathcal{N}_x$ ,
- (N2) The intersection of any two neighborhoods of  $x$  is another neighborhood of  $x$ :  $N, M \in \mathcal{N}_x \Rightarrow N \cap M \in \mathcal{N}_x$ ,
- (N3) Every superset of any neighborhood of  $x$  is a neighborhood of  $x$ :  $(M \in \mathcal{N}_x) \wedge (M \subseteq N) \Rightarrow N \in \mathcal{N}_x$ ,

that characterize  $\mathcal{N}_x$  completely are a direct consequence of the definitions (15), (16) that may also be stated as

(N0) Any neighborhood  $N \in \mathcal{N}_x$  contains another neighborhood  $U$  of  $x$  that is a *neighborhood of each of its points*:  $((\forall N \in \mathcal{N}_x)(\exists U \in \mathcal{N}_x)(U \subseteq N)) : (\forall y \in U \Rightarrow U \in \mathcal{N}_y)$ .

Property (N0) infact serves as the defining characteristic of an open set, and  $U$  can be identified with the largest open set  $\text{Int}(N)$  contained in  $N$ ; hence a set  $G$  in a topological space is open iff it is a neighborhood of each of its points. Accordingly if  $\mathcal{N}_x$  is a given class of subsets of  $X$  associated with each  $x \in X$  satisfying (N1)–(N3), then (N0) defines the special class of neighborhoods  $G$

$$\mathcal{U} = \{G \in \mathcal{N}_x : x \in B \subseteq G \text{ for all } x \in G \text{ and some basic nbd } B \in \mathcal{N}_x\} \tag{18}$$

as the unique topology on  $X$  that contains a basic neighborhood of each of its points, for which the

neighborhood system at  $x$  coincides exactly with the assigned collection  $\mathcal{N}_x$ ; compare with Definition A.1.1. Neighborhoods in topological spaces are a generalization of the familiar notion of distances of metric spaces that quantifies “closeness” of points of  $X$ .

A *neighborhood of a non-empty subset  $A$*  of  $X$  that will be needed later on is defined in a similar manner:  $N$  is a neighborhood of  $A$  iff  $A \subseteq \text{Int}(N)$ , that is  $A \subseteq U \subseteq N$ ; thus the neighborhood system at  $A$  is given by  $\mathcal{N}_A = \bigcap_{a \in A} \mathcal{N}_a := \{G \subseteq X : G \in \mathcal{N}_a \text{ for every } a \in A\}$  is the class of common neighborhoods of each point of  $A$ .

Some examples of neighborhood systems at a point  $x$  in  $X$  are the following:

- (1) In an indiscrete space  $(X, \mathcal{U})$ ,  $X$  is the only neighborhood of every point of the space; in a discrete space any set containing  $x$  is a neighborhood of the point.
- (2) In an infinite cofinite (or uncountable cocountable) space, every neighborhood of a point is an open neighborhood of that point.
- (3) In the topology of saturated sets under the equivalence relation  $\sim$ , the neighborhood system at  $x$  consists of all supersets of the equivalence class  $[x]_{\sim}$ .
- (4) Let  $x \in X$ . In the  $x$ -inclusion topology,  $\mathcal{N}_x$  consists of all the non-empty open sets of  $X$  which are the supersets of  $\{x\}$ . For a point  $y \neq x$  of  $X$ ,  $\mathcal{N}_y$  are the supersets of  $\{x, y\}$ .

For any given class  $_{\text{T}}\mathcal{S}$  of subsets of  $X$ , a unique topology  $\mathcal{U}(_{\text{T}}\mathcal{S})$  can always be constructed on  $X$  by taking all *finite intersections*  $_{\text{T}}\mathcal{S}_{\wedge}$  of members of  $\mathcal{S}$  followed by *arbitrary unions*  $_{\text{T}}\mathcal{S}_{\vee}$  of these finite intersections.  $\mathcal{U}(_{\text{T}}\mathcal{S}) := _{\text{T}}\mathcal{S}_{\wedge\vee}$  is the smallest topology on  $X$  that contains  $_{\text{T}}\mathcal{S}$  and is said to be *generated by*  $_{\text{T}}\mathcal{S}$ . For a given topology  $\mathcal{U}$  on  $X$  satisfying  $\mathcal{U} = \mathcal{U}(_{\text{T}}\mathcal{S})$ ,  $_{\text{T}}\mathcal{S}$  is a *subbasis*, and  $_{\text{T}}\mathcal{S}_{\wedge} := _{\text{T}}\mathcal{B}$  a *basis, for the topology  $\mathcal{U}$* ; for more on topological basis, see Appendix A.1. The topology generated by a subbase essentially builds not from the collection  $_{\text{T}}\mathcal{S}$  itself but from the finite intersections  $_{\text{T}}\mathcal{S}_{\wedge}$  of its subsets; in comparison the base generates a topology directly from a collection  $_{\text{T}}\mathcal{S}$  of subsets by forming their unions. Thus whereas *any* class of subsets can be used as a subbasis, a given collection must meet certain qualifications to pass the test of a base for a topology: these and related topics are covered in Appendix A.1. Different subbases, therefore, can be used to generate different topologies on the same set  $X$  as the following examples for the case of



$X = \mathbb{R}$  demonstrates; here  $(a, b)$ ,  $[a, b)$ ,  $(a, b]$  and  $[a, b]$ , for  $a \leq b \in \mathbb{R}$ , are the usual open-closed intervals in  $\mathbb{R}$ .<sup>13</sup> The subbases  ${}_T\mathcal{S}_1 = \{(a, \infty), (-\infty, b)\}$ ,  ${}_T\mathcal{S}_2 = \{[a, \infty), (-\infty, b]\}$ ,  ${}_T\mathcal{S}_3 = \{(a, \infty), (-\infty, b]\}$  and  ${}_T\mathcal{S}_4 = \{[a, \infty), (-\infty, b]\}$  give the respective bases  ${}_T\mathcal{B}_1 = \{(a, b)\}$ ,  ${}_T\mathcal{B}_2 = \{[a, b)\}$ ,  ${}_T\mathcal{B}_3 = \{(a, b]\}$  and  ${}_T\mathcal{B}_4 = \{[a, b]\}$ ,  $a \leq b \in \mathbb{R}$ , leading to the *standard (usual)*, *lower limit (Sorgenfrey)*, *upper limit*, and *discrete* (take  $a = b$ ) topologies on  $\mathbb{R}$ . Bases of the type  $(a, \infty)$  and  $(-\infty, b)$  provide the *right* and *left ray* topologies on  $\mathbb{R}$ .

*This feasibility of generating different topologies on a set can be of great practical significance because open sets determine convergence characteristics of nets and continuity characteristics of functions, thereby making it possible for nature to play around with the structure of its working space in its kitchen to its best possible advantage.*<sup>14</sup>

Here are a few essential concepts and terminology for topological spaces.

**Definition 2.2** (Boundary, Closure, Interior). The boundary of  $A$  in  $X$  is the set of points  $x \in X$  such that every neighborhood  $N$  of  $x$  intersects both  $A$  and  $X - A$ :

$$\text{Bdy}(A) \stackrel{\text{def}}{=} \{x \in X : (\forall N \in \mathcal{N}_x)((N \cap A \neq \emptyset \wedge (N \cap (X - A) \neq \emptyset))\} \quad (19)$$

where  $\mathcal{N}_x$  is the neighborhood system of Eq. (17) at  $x$ .

The closure of  $A$  is the set of all points  $x \in X$  such that each neighborhood of  $x$  contains at least one point of  $A$  *that may be  $x$  itself*. Thus the set

$$\text{Cl}(A) \stackrel{\text{def}}{=} \{x \in X : (\forall N \in \mathcal{N}_x)(N \cap A \neq \emptyset)\} \quad (20)$$

of all points in  $X$  adherent to  $A$  is given by the union of  $A$  with its boundary.

The interior of  $A$

$$\text{Int}(A) \stackrel{\text{def}}{=} \{x \in X : (\exists N \in \mathcal{N}_x)(N \subseteq A)\} \quad (21)$$

consisting of those points of  $X$  that are in  $A$  but not in its boundary,  $\text{Int}(A) = A - \text{Bdy}(A)$ , is the largest open subset of  $X$  that is contained in  $A$ . Hence it follows that  $\text{Int}(\text{Bdy}(A)) = \emptyset$ , the boundary of  $A$  is the intersection of the closures of  $A$  and  $X - A$ , and a subset  $N$  of  $X$  is a neighborhood of  $x$  iff  $x \in \text{Int}(N)$ .

The three subsets  $\text{Int}(A)$ ,  $\text{Bdy}(A)$  and *exterior* of  $A$  defined as  $\text{Ext}(A) := \text{Int}(X - A) = X - \text{Cl}(A)$ , are pairwise disjoint and have the full space  $X$  as their union.

**Definition 2.3** (Derived and Isolated sets). Let  $A$  be a subset of  $X$ . A point  $x \in X$  (which may or may not be a point of  $A$ ) is a cluster point of  $A$  if every neighborhood  $N \in \mathcal{N}_x$  contains at least one point of  $A$  *different from  $x$* . The derived set of  $A$

$$\text{Der}(A) \stackrel{\text{def}}{=} \left\{x \in X : (\forall N \in \mathcal{N}_x)\left(N \cap (A - \{x\}) \neq \emptyset\right)\right\} \quad (22)$$

is the set of all cluster points of  $A$ . The complement of  $\text{Der}(A)$  in  $A$

$$\text{Iso}(A) \stackrel{\text{def}}{=} A - \text{Der}(A) = \text{Cl}(A) - \text{Der}(A) \quad (23)$$

are the isolated points of  $A$  to which no proper sequence in  $A$  converges, that is there exists a neighborhood of any such point that contains no other point of  $A$  so that the only sequence that converges to  $a \in \text{Iso}(A)$  is the constant sequence  $(a, a, a, \dots)$ .

Clearly,

$$\begin{aligned} \text{Cl}(A) &= A \cup \text{Der}(A) = A \cup \text{Bdy}(A) \\ &= \text{Iso}(A) \cup \text{Der}(A) = \text{Int}(A) \cup \text{Bdy}(A) \end{aligned}$$

with the last two being disjoint unions, and  $A$  is closed iff  $A$  contains all its cluster points,  $\text{Der}(A) \subseteq A$ , iff  $A$  contains its closure. Hence

$$\begin{aligned} A = \text{Cl}(A) &\Leftrightarrow \text{Cl}(A) \\ &= \{x \in A : ((\exists N \in \mathcal{N}_x)(N \subseteq A)) \\ &\quad \vee ((\forall N \in \mathcal{N}_x)(N \cap (X - A) \neq \emptyset))\}. \end{aligned}$$

<sup>13</sup>By definition, an interval  $I$  in a totally ordered set  $X$  is a subset of  $X$  with the property

$$(x_1, x_2 \in I) \wedge (x_3 \in X : x_1 \prec x_3 \prec x_2) \Rightarrow x_3 \in I$$

so that any element of  $X$  lying between two elements of  $I$  also belongs to  $I$ .

<sup>14</sup>Although we do not pursue this point of view here, it is nonetheless tempting to speculate that the answer to the question “Why does the entropy of an isolated system increase?” may be found by exploiting this line of reasoning that seeks to explain the increase in terms of a visible component associated with the usual topology as against a different latent workplace topology that governs the dynamics of nature.

Comparison of Eqs. (19) and (22) also makes it clear that  $\text{Bdy}(A) \subseteq \text{Der}(A)$ . The special case of  $A = \text{Iso}(A)$  with  $\text{Der}(A) \subseteq X - A$  is important enough to deserve a special mention:

**Definition 2.4** (Donor set). A proper, nonempty subset  $A$  of  $X$  such that  $\text{Iso}(A) = A$  with  $\text{Der}(A) \subseteq X - A$  will be called self-isolated or donor. Thus sequences eventually in a donor set converges only in its complement; this is, the opposite of the characteristic of a closed set where all converging sequences eventually in the set must necessarily converge in it. A closed-donor set with a closed neighbor has no derived or boundary sets, and will be said to be isolated in  $X$ .

**Example 2.5.** In an isolated set sequences converge, if they have to, simultaneously in the complement (because it is donor) and in it (because it is closed). Convergent sequences in such a set can only be constant sequences. Physically, if we consider adherents to be contributions made by the dynamics of the corresponding sequences, then an isolated set is secluded from its neighbor in the sense that it neither receives any contributions from its surroundings, nor does it give away any. In this light and terminology, a closed set is a *selfish* set (recall that a set  $A$  is closed in  $X$  iff every convergent net of  $X$  that is eventually in  $A$  converges in  $A$ ; conversely a set is open in  $X$  iff the only nets that converge in  $A$  are eventually in it), whereas a set with a derived set that intersects itself and its complement may be considered to be *neutral*. Appendix A.3 shows the various possibilities for the derived set and boundary of a subset  $A$  of  $X$ .

Some useful properties of these concepts for a subset  $A$  of a topological space  $X$  are the following.

- (a)  $\text{Bdy}_X(X) = \emptyset$ ,
  - (b)  $\text{Bdy}(A) = \text{Cl}(A) \cap \text{Cl}(X - A)$ ,
  - (c)  $\text{Int}(A) = X - \text{Cl}(X - A) = A - \text{Bdy}(A) = \text{Cl}(A) - \text{Bdy}(A)$ ,
  - (d)  $\text{Int}(A) \cap \text{Bdy}(A) = \emptyset$ ,
  - (e)  $X = \text{Int}(A) \cup \text{Bdy}(A) \cup \text{Int}(X - A)$ ,
  - (f)  $\text{Int}(A) = \bigcup \{G \subseteq X : G$   
is an open set of  $X$  contained in  $A\}$
- (24)

$$(g) \quad \text{Cl}(A) = \bigcap \{F \subseteq X : F \text{ is a closed set of } X \text{ containing } A\} \quad (25)$$

A straightforward consequence of property (b) is that the boundary of any subset  $A$  of a topological space  $X$  is closed in  $X$ ; this significant result may also be demonstrated as follows. If  $x \in X$  is not in the boundary of  $A$  there is some neighborhood  $N$  of  $x$  that does not intersect both  $A$  and  $X - A$ . For each point  $y \in N$ ,  $N$  is a neighborhood of that point that does not meet  $A$  and  $X - A$  simultaneously so that  $N$  is contained wholly in  $X - \text{Bdy}(A)$ . We may now take  $N$  to be open without any loss of generality implying thereby that  $X - \text{Bdy}(A)$  is an open set of  $X$  from which it follows that  $\text{Bdy}(A)$  is closed in  $X$ .

Further material on topological spaces relevant to our work can be found in Appendix A.3.

### End Tutorial 4

---

Working in a general topological space, we now recall the solution of an ill-posed problem  $f(x) = y$  [Sengupta, 1997] that leads to a multifunctional inverse  $f^-$  through the generalized inverse  $G$ . Let  $f: (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  be a (nonlinear) function between two topological space  $(X, \mathcal{U})$  and  $(Y, \mathcal{V})$  that is neither one-one or onto. Since  $f$  is not one-one,  $X$  can be partitioned into disjoint equivalence classes with respect to the equivalence relation  $x_1 \sim x_2 \Leftrightarrow f(x_1) = f(x_2)$ . Picking a representative member from each of the classes (this is possible by the Axiom of Choice; see the following Tutorial) produces a *basic set*  $X_B$  of  $X$ ; it is basic as it corresponds to the row space in the linear matrix example which is all that is needed for taking an inverse.  $X_B$  is the counterpart of the quotient set  $X/\sim$  of Sec. 1, with the important difference that whereas the points of the quotient set are the equivalence classes of  $X$ ,  $X_B$  is a subset of  $X$  with each of the classes contributing a point to  $X_B$ . It then follows that  $f_B: X_B \rightarrow f(X)$  is the bijective restriction  $a|_{\text{row}(A)}$  that reduces the original ill-posed problem to a well-posed one with  $X_B$  and  $f(X)$  corresponding respectively to the row and column spaces of  $A$ , and  $f_B^{-1}: f(X) \rightarrow X_B$  is the basic inverse from which the multiinverse  $f^-$  is obtained through  $G$ , which in turn corresponds to the Moore–Penrose inverse  $G_{\text{MP}}$ . The topological considerations (obviously not for inner product spaces

that applies to the Moore–Penrose inverse) needed to complete the solution are discussed below and in Appendix A.1.

---

### Begin Tutorial 5: Axiom of Choice and Zorn’s Lemma

Since some of our basic arguments depend on it, this Tutorial contains a short description of the Axiom of Choice that has been described as “one of the most important, and at the same time one of the most controversial, principles of mathematics”. What this axiom states is this: For any set  $X$  there exists a function  $f_C : \mathcal{P}_0(X) \rightarrow X$  such that  $f_C(A_\alpha) \in A_\alpha$  for every non-empty subset  $A_\alpha$  of  $X$ ; here  $\mathcal{P}_0(X)$  is the class of all subsets of  $X$  except  $\emptyset$ . Thus, if  $X = \{x_1, x_2, x_3\}$  is a three element set, a possible choice function is given by

$$f_C(\{x_1, x_2, x_3\}) = x_3, \quad f_C(\{x_1, x_2\}) = x_1,$$

$$f_C(\{x_2, x_3\}) = x_3, \quad f_C(\{x_3, x_1\}) = x_3,$$

$$f_C(\{x_1\}) = x_1, \quad f_C(\{x_2\}) = x_2, \quad f_C(\{x_3\}) = x_3.$$

It must be appreciated that the axiom is only an existence result that asserts *every set* to have a choice function, even when nobody knows how to construct one in a specific case. Thus, for example, how does one pick out the isolated irrationals  $\sqrt{2}$  or  $\pi$  from the uncountable reals? There is no doubt that they do exist, for we can construct a right-angled triangle with sides of length 1 or a circle of radius 1. The axiom tells us that these choices are possible even though we do not know how exactly to do it; all that can be stated with confidence is that we can actually pick up rationals arbitrarily close to these irrationals.

The axiom of choice is essentially meaningful when  $X$  is infinite as illustrated in the last two examples. This is so because even when  $X$  is denumerable, it would be physically impossible to make an infinite number of selections either all at a time or sequentially: the Axiom of Choice nevertheless tells us that this is possible. The real strength and utility of the Axiom however is when  $X$  and some or all of its subsets are uncountable as in the case

of the choice of the *single element*  $\pi$  from the reals. To see this more closely in the context of maps that we are concerned with, let  $f : X \rightarrow Y$  be a non-injective, onto map. To construct a functional right inverse  $f_r : Y \rightarrow X$  of  $f$ , we must choose, for each  $y \in Y$  one *representative* element  $x_{\text{rep}}$  from the set  $f^{-}(y)$  and define  $f_r(y)$  to be that element according to  $f \circ f_r(y) = f(x_{\text{rep}}) = y$ . If there is no preferred or natural way to make this choice, the axiom of choice allows us to make an arbitrary selection from the infinitely many that may be possible from  $f^{-}(y)$ . When a natural choice is indeed available, as for example in the case of the initial value problem  $y'(x) = x; y(0) = \alpha_0$  on  $[0, a]$ , the definite solution  $\alpha_0 + x^2/2$  may be selected from the infinitely many  $\int_0^x x'dx' = \alpha + x^2/2, 0 \leq x \leq a$  that are permissible, and the axiom of choice sanctions this selection. In addition, each  $y \in Y$  gives rise to the family of solution sets  $A_y = \{f^{-}(y) : y \in Y\}$  and the real power of the axiom is its assertion that it is possible to make a choice  $f_C(A_y) \in A_y$  on every  $A_y$  simultaneously; this permits the choice on every  $A_y$  of the collection to be made at the same time.

### Pause Tutorial 5

---

Figure shows our formulation and solution [Sengupta, 1997] of the inverse ill-posed problem  $f(x) = y$ . In sub-diagram  $X - X_B - f(X)$ , the surjection  $p : X \rightarrow X_B$  is the counterpart of the quotient map  $Q$  of Fig. 2 that is known in the present context as the *identification* of  $X$  with  $X_B$  (as it *identifies* each saturated subset of  $X$  with its representative point in  $X_B$ ), with the space  $(X_B, \text{FT}\{\mathcal{U}; p\})$  carrying the *identification topology*  $\text{FT}\{\mathcal{U}; p\}$  being known as an *identification space*. By sub-diagram  $Y - X_B - f(X)$ , the image  $f(X)$  of  $f$  gets the *subspace topology*<sup>15</sup>  $\text{IT}\{j; \mathcal{V}\}$  from  $(Y, \mathcal{V})$  by the inclusion  $j : f(X) \rightarrow Y$  when its open sets are generated as, and only as,  $j^{-1}(V) = V \cap f(X)$  for  $V \in \mathcal{V}$ . Furthermore if the bijection  $f_B$  connecting  $X_B$  and  $f(X)$  (which therefore acts as a 1 : 1 correspondence between their points, implying that these sets are set-theoretically identical

---

<sup>15</sup>In a subspace  $A$  of  $X$ , a subset  $U_A$  of  $A$  is open iff  $U_A = A \cap U$  for some open set  $U$  of  $X$ . The notion of subspace topology can be formalized with the help of the inclusion map  $i : A \rightarrow (X, \mathcal{U})$  that puts every point of  $A$  back to where it came from, thus

$$\begin{aligned} \mathcal{U}_A &= \{U_A = A \cap U : U \in \mathcal{U}\} \\ &= \{i^{-}(U) : U \in \mathcal{U}\}. \end{aligned}$$

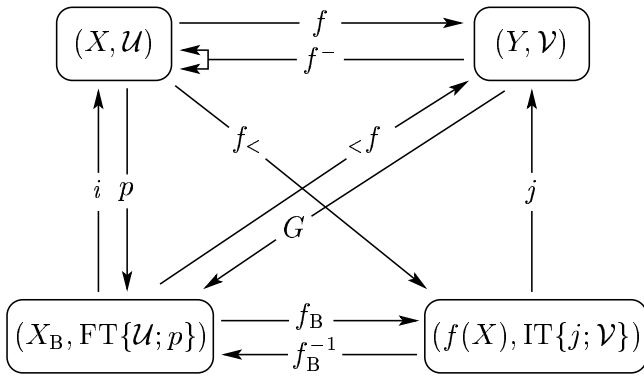


Fig. 5. Solution of ill-posed problem  $f(x) = y, f : X \rightarrow Y$ .  $G : Y \rightarrow X_B$ , a generalized inverse of  $f$  because of  $fGf = f$  and  $GfG = G$  which follows from the commutativity of the diagrams, is a functional selection of the multi-inverse  $f^- : (Y, \mathcal{V}) \multimap (X, \mathcal{U})$   $<f$  and  $f_<$  are the injective and surjective restrictions of  $f$ ; these will be topologically denoted by their generic notations  $e$  and  $q$ , respectively.

except for their names) is image continuous, then by Theorem A.2.1 of Appendix 2, so is the *association*  $q = f_B \circ p : X \rightarrow f(X)$  that associates saturated sets of  $X$  with elements of  $f(X)$ ; this makes  $f(X)$  look like an identification space of  $X$  by assigning to it the topology  $\text{FT}\{\mathcal{U}; q\}$ . On the other hand if  $f_B$  happens to be preimage continuous, then  $X_B$  acquires, by Theorem A.2.2, the initial topology  $\text{IT}\{e; \mathcal{V}\}$  by the *embedding*  $e : X_B \rightarrow Y$  that embeds  $X_B$  into  $Y$  through  $j \circ f_B$ , making it look like a subspace of  $Y$ .<sup>16</sup> In this dual situation,  $f_B$  has the highly interesting topological property of being simultaneously image and preimage continuous when the open sets of  $X_B$  and  $f(X)$  — which are simply the  $f_B^{-1}$ -images of the open sets of  $f(X)$  which, in turn, are the  $f_B$ -images of these saturated open sets — can be considered to have been generated by  $f_B$ , and are respectively the smallest and largest collection of subsets of  $X$  and  $Y$  that makes  $f_B$  *initial-final continuous* [Sengupta, 1997]. A bijective ininal function such as  $f_B$  is known as a *homeomorphism* and ininality for functions that are neither 1 : 1 nor onto is a generalization of homeomorphism for bijections; refer Eqs. (A.47) and (A.48) for a set-theoretic formulation of this distinction. A homeomorphism  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  renders the homeomorphic spaces  $(X, \mathcal{U})$  and  $(Y, \mathcal{V})$  topologically

indistinguishable which may be considered to be identical in as far as their topological properties are concerned.

*Remark.* It may be of some interest here to speculate on the significance of *ininality* in our work. Physically, a map  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  between two spaces can be taken to represent an interaction between them and the algebraic and topological characters of  $f$  determine the nature of this interaction. A simple bijection merely sets up a correspondence, that is an interaction, between every member of  $X$  with some member  $Y$ , whereas a continuous map establishes the correspondence among the special category of “open” sets. Open sets, as we see in Appendix A.1, are the basic ingredients in the theory of convergence of sequences, nets and filters, and the characterization of open sets in terms of convergence, namely that *a set  $G$  in  $X$  is open in it if every net or sequence that converges in  $X$  to a point in  $G$  is eventually in  $G$* , see Appendix A.1, may be interpreted to mean that such sets represent groupings of elements that require membership of the group before permitting an element to belong to it; an open set unlike its complement the closed or *selfish* set, however, does not forbid a net that has been eventually in it to settle down in its selfish neighbor, who nonetheless will never allow such a situation to develop in its own territory. An ininal map forces these well-defined and definite groups in  $(X, \mathcal{U})$  and  $(Y, \mathcal{V})$  to interact with each other through  $f$ ; this is not possible with simple continuity as there may be open sets in  $X$  that are not derived from those of  $Y$  and non-open sets in  $Y$  whose inverse images are open in  $X$ . *It is our hypothesis that the driving force behind the evolution of a system represented by the input-output relation  $f(x) = y$  is the attainment of the ininal triple state  $(X, f, Y)$  for the system.* A preliminary analysis of this hypothesis is to be found in Sec. 4.2.

For ininality of the interaction, it is therefore necessary to have

$$\begin{aligned} \text{FT}\{\mathcal{U}; f_<\} &= \text{IT}\{j; \mathcal{V}\} \\ \text{IT}\{<f; \mathcal{V}\} &= \text{FT}\{\mathcal{U}; p\}; \end{aligned} \tag{26}$$

in what follows we will refer to the injective and surjective restrictions of  $f$  by their generic topological symbols of embedding  $e$  and association  $q$ , respectively. What are the topological characteristics of  $f$

<sup>16</sup>A surjective function is an *association* iff it is image continuous and an injective function is an *embedding* iff it is preimage continuous.

in order that the requirements of Eq. (26) be met? From Appendix A.1, it should be clear by superposing the two parts of Fig. 21 over each other that given  $q : (X, \mathcal{U}) \rightarrow (f(X), \text{FT}\{\mathcal{U}; q\})$  in the first of these equations,  $\text{IT}\{j; \mathcal{V}\}$  will equal  $\text{FT}\{\mathcal{U}; q\}$  iff  $j$  is an ininal open inclusion and  $Y$  receives  $\text{FT}\{\mathcal{U}; f\}$ . In a similar manner, preimage continuity of  $e$  requires  $p$  to be open ininal and  $f$  to be preimage continuous if the second of Eq. (26) is to be satisfied. Thus under the restrictions imposed by Eq. (26), the interaction  $f$  between  $X$  and  $Y$  must be such as to give  $X$  the smallest possible topology of  $f$ -saturated sets and  $Y$  the largest possible topology of images of all these sets:  $f$ , under these conditions, is an ininal transformation. Observe that a direct application of parts (b) of Theorems A.2.1 and A.2.2 to Fig. implies that Eq. (26) is satisfied iff  $f_B$  is ininal, that is iff it is a homeomorphism. Ininality of  $f$  is simply a reflection of this as it is neither 1 : 1 nor onto.

The  $f$ - and  $p$ -images of each saturated set of  $X$  are singletons in  $Y$  (these saturated sets in  $X$  arose, in the first place, as  $f^-(\{y\})$  for  $y \in Y$ ) and in  $X_B$ , respectively. This permits the embedding  $e = j \circ f_B$  to give  $X_B$  the character of a virtual subspace of  $Y$  just as  $i$  makes  $f(X)$  a real subspace. Hence the inverse images  $p^-(x_r) = f^-(e(x_r))$  with  $x_r \in X_B$ , and  $q^-(y) = f^-(i(y))$  with  $y = f_B(x_r) \in f(X)$  are the same, and are just the corresponding  $f^-$  images via the injections  $e$  and  $i$ , respectively.  $G$ , a left inverse of  $e$ , is a generalized inverse of  $f$ .  $G$  is a generalized inverse because the two set-theoretic defining requirements of  $fGf = f$  and  $GfG = G$  for the generalized inverse are satisfied, as Fig. shows, in the following forms

$$jf_BGf = f \quad Gj_BG = G.$$

In fact the commutativity embodied in these equalities is self evident from the fact that  $e = if_B$  is a left inverse of  $G$ , that is  $eG = \mathbf{1}_Y$ . On putting back  $X_B$  into  $X$  by identifying each point of  $X_B$  with the set it came from yields the required set-valued inverse  $f^-$ , and  $G$  may be viewed as a functional selection of the multiinverse  $f^-$ .

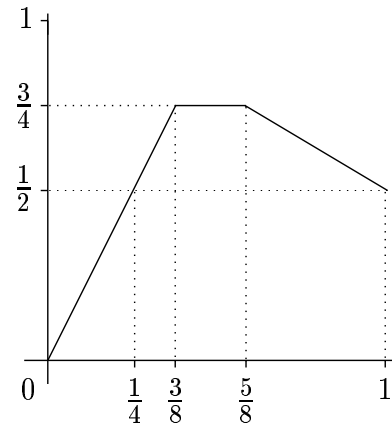


Fig. 6. The function  $f(x) = \begin{cases} 2x, & 0 \leq x < 3/8 \\ 3/4, & 3/8 \leq x \leq 5/8 \\ 7/6 - 2x/3, & 5/8 < x \leq 1. \end{cases}$

An *injective branch* of a function  $f$  in this work refers to the restrictions  $f_B$  and its associated inverse  $f_B^{-1}$ .

The following example of an inverse ill-posed problem will be useful in fixing the notations introduced above. Let  $f$  on  $[0, 1]$  be the function shown below.

Then  $f(x) = y$  is well-posed for  $[0, 1/4)$ , and ill-posed in  $[1/4, 1]$ . There are two injective branches of  $f$  in  $\{[1/4, 3/8) \cup (5/8, 1]\}$ , and  $f$  is constant ill-posed in  $[3/8, 5/8]$ . Hence the basic component  $f_B$  of  $f$  can be taken to be  $f_B(x) = 2x$  for  $x \in [0, 3/8)$  having the inverse  $f_B^{-1}(y) = x/2$  with  $y \in [0, 3/4]$ . The generalized inverse is obtained by taking  $[0, 3/4]$  as a subspace of  $[0, 1]$ , while the multiinverse  $f^-$  follows by associating with every point of the basic domain  $[0, 1]_B = [0, 3/8]$ , the respective equivalent points  $[3/8]_f = [3/8, 5/8]$  and  $[x]_f = \{x, 7/4 - 3x\}$  for  $x \in [1/4, 3/8)$ . Thus the inverses  $G$  and  $f^-$  of  $f$  are<sup>17</sup>

$$G(y) = \begin{cases} \frac{y}{2}, & y \in \left[0, \frac{3}{4}\right] \\ 0, & y \in \left(\frac{3}{4}, 1\right] \end{cases},$$

<sup>17</sup>If  $y \notin \mathcal{R}(f)$  then  $f^-(\{y\}) := \emptyset$  which is true for any subset of  $Y - \mathcal{R}(f)$ . However from the set-theoretic definition of natural numbers that requires  $0 := \emptyset, 1 = \{0\}, 2 = \{0, 1\}$  to be defined recursively, it follows that  $f^-(y)$  can be identified with 0 whenever  $y$  is not in the domain of  $f^-$ . Formally, the successor set  $A^+ = A \cup \{A\}$  of  $A$  can be used to write  $0 := \emptyset, 1 = 0^+ = 0 \cup \{0\}, 2 = 1^+ = 1 \cup \{1\} = \{0\} \cup \{1\}, 3 = 2^+ = 2 \cup \{2\} = \{0\} \cup \{1\} \cup \{2\}$ , etc. Then the set of natural numbers  $\mathbb{N}$  is defined to be the intersection of all the successor sets, where a successor set  $\mathcal{S}$  is any set that contains  $\emptyset$  and  $A^+$  whenever  $A$  belongs to  $\mathcal{S}$ . Observe how in the successor notation, countable union of singleton integers recursively define the corresponding sum of integers.

$$f^{-}(y) = \begin{cases} \frac{y}{2}, & y \in \left[0, \frac{1}{2}\right) \\ \left\{\frac{y}{2}, \frac{7}{4} - \frac{3y}{2}\right\}, & y \in \left[\frac{1}{2}, \frac{3}{4}\right) \\ \left[\frac{3}{8}, \frac{5}{8}\right], & y = \frac{3}{4} \\ 0, & y \in \left(\frac{3}{4}, 1\right], \end{cases}$$

which shows that  $f^{-}$  is multivalued. In order to avoid cumbersome notations, an injective branch of  $f$  will always refer to a representative basic branch  $f_B$ , and its “inverse” will mean either  $f_B^{-1}$  or  $G$ .

**Example 2.3** (Revisited). The row reduced echelon form of the augmented matrix  $(A|b)$  of Example 2.3 is

$$(A|b) \rightarrow \begin{pmatrix} 1 & -3 & 0 & \frac{3}{2} & \frac{1}{2} & \frac{5b_1}{2} - \frac{b_2}{2} \\ 0 & 0 & 1 & -\frac{1}{4} & \frac{3}{4} & -\frac{3b_1}{4} + \frac{b_2}{4} \\ 0 & 0 & 0 & 0 & 0 & -2b_1 + b_3 \\ 0 & 0 & 0 & 0 & 0 & b_1 - b_2 + b_4 \end{pmatrix} \quad (27)$$

The multifunctional solution  $x = A^{-}b$ , with  $b$  any element of  $Y = \mathbb{R}^4$  not necessarily in the image of  $a$ , is

$$x = A^{-}b = Gb + x_2 \begin{pmatrix} 3 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix} + x_4 \begin{pmatrix} -\frac{3}{2} \\ 0 \\ \frac{1}{4} \\ 1 \\ 0 \end{pmatrix} + x_5 \begin{pmatrix} -\frac{1}{2} \\ 0 \\ -\frac{3}{4} \\ 0 \\ 1 \end{pmatrix},$$

with its multifunctional character arising from the arbitrariness of the coefficients  $x_2, x_4$  and  $x_5$ . The generalized inverse

$$G = \begin{pmatrix} \frac{5}{2} & -\frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -\frac{3}{4} & \frac{1}{4} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} : Y \rightarrow X_B \quad (28)$$

is the unique matrix representation of the functional inverse  $a_B^{-1} : a(\mathbb{R}^5) \rightarrow X_B$  extended to  $Y$  defined according to<sup>18</sup>

$$g(b) = \begin{cases} a_B^{-1}(b), & \text{if } b \in \mathcal{R}(a) \\ 0, & \text{if } b \in Y - \mathcal{R}(a), \end{cases} \quad (29)$$

that bears comparison with the basic inverse

$$A_B^{-1}(b^*) = \begin{pmatrix} \frac{5}{2} & -\frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ -\frac{3}{4} & \frac{1}{4} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \times \begin{pmatrix} b_1 \\ b_2 \\ 2b_1 \\ b_2 - b_1 \end{pmatrix} : a(\mathbb{R}^5) \rightarrow X_B$$

between the two-dimensional column and row spaces of  $A$  which is responsible for the particular solution of  $Ax = b$ . Thus  $G$  is simply  $A_B^{-1}$  acting on its domain  $a(X)$  considered a subspace of  $Y$ , suitably extended to the whole of  $Y$ . That it is indeed a generalized inverse is readily seen through the matrix multiplications  $GAG$  and  $AGA$  that can be verified to reproduce  $G$  and  $A$ , respectively. Comparison of Eqs. (12) and (29) shows that the Moore–Penrose inverse differs from ours through the geometrical constraints imposed in its definition, Eqs. (13). Of course, this results in a more complex inverse (14) as compared to our very simple (28); nevertheless it is true that both the inverses satisfy

$$\begin{aligned} E((E(G_{MP}))^T) &= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix} \\ &= E((E(G))^T) \end{aligned}$$

where  $E(A)$  is the row-reduced echelon form of  $A$ . The canonical simplicity of Eq. (28) as compared to Eq. (14) is a general feature that suggests a more natural choice of bases by the map  $a$  than the orthogonal set imposed by Moore and Penrose. This is to be expected since the MP inverse, governed by Eq. (13), is a subset of our less restricted inverse

<sup>18</sup>See footnote 17 for a justification of the definition when  $b$  is not in  $\mathcal{R}(a)$ .

described by only the first two of (13); more specifically the difference is made clear in Fig. 4(a) which shows that for any  $b \notin \mathcal{R}(A)$ , only  $G_{\text{MP}}(b_{\perp}) = 0$  as compared to  $G(b) = 0$ . This seems to imply that introducing extraneous topological considerations into the purely set-theoretic inversion process may not be a recommended way of inverting, and the simple bases comprising the row and null spaces of  $A$  and  $A^{\text{T}}$  — that are mutually orthogonal just as those of the Moore–Penrose — are a better choice for the particular problem  $Ax = b$  than the general orthonormal bases that the MP inverse introduces. These “good” bases, with respect to which the generalized inverse  $G$  has a considerably simpler representation, are obtained in a straightforward manner from the row-reduced forms of  $A$  and  $A^{\text{T}}$ . These bases are

- (a) The column space of  $A$  is spanned by the columns  $(1, 3, 2, 2)^{\text{T}}$  and  $(1, 5, 2, 4)^{\text{T}}$  of  $A$  that correspond to the basic columns containing the leading 1’s in the row-reduced form of  $A$ ,
- (b) The null space of  $A^{\text{T}}$  is spanned by the solutions  $(-2, 0, 1, 0)^{\text{T}}$  and  $(1, -1, 0, 1)^{\text{T}}$  of the equation  $A^{\text{T}}b = 0$ ,
- (c) The row space of  $A$  is spanned by the rows  $(1, -3, 2, 1, 2)$  and  $(3, -9, 10, 2, 9)$  of  $A$  corresponding to the non-zero rows in the row-reduced form of  $A$ ,
- (d) The null space of  $A$  is spanned by the solutions  $(3, 1, 0, 0, 0)$ ,  $(-6, 0, 1, 4, 0)$ , and  $(-2, 0, -3, 0, 4)$  of the equation  $Ax = 0$ .

The main differences between the natural “good” bases and the MP-bases that are responsible for the difference in the form of inverses, is that the latter have the additional restrictions of being orthogonal to each other (recall the orthogonality property of the  $Q$ -matrices), and the more severe of basis vectors mapping onto basis vectors according to  $Ax_i = \sigma_i b_i$ ,  $i = 1, \dots, r$ , where the  $\{x_i\}_{i=1}^n$  and  $\{b_j\}_{j=1}^m$  are the eigenvectors of  $A^{\text{T}}A$  and  $AA^{\text{T}}$  respectively and  $(\sigma_i)_{i=1}^r$  are the positive square roots of the non-zero eigenvalues of  $A^{\text{T}}A$  (or of  $AA^{\text{T}}$ ), with  $r$  denoting the dimension of the row or column space. This is considered as a serious restriction as the linear combination of the basis  $\{b_j\}$  that  $Ax_i$  should otherwise have been equal to, allows a greater flexibility in the matrix representation of the inverse that shows up in the structure of  $G$ . These are, in fact, quite general considerations in the matrix representation of linear operators; thus

the basis that diagonalizes an  $n \times n$  matrix (when this is possible) is not the standard “diagonal” orthonormal basis of  $\mathbb{R}^n$ , but a problem-dependent, less canonical, basis consisting of the  $n$  eigenvectors of the matrix. The 0-rows of the inverse of Eq. (28) result from the three-dimensional null-space variables  $x_2, x_4$  and  $x_5$ , while the 0-columns come from the two-dimensional image-space dependency of  $b_3, b_4$  on  $b_1$  and  $b_2$ , that is from the last two zero rows of the reduced echelon form (27) of the augmented matrix.

We will return to this theme of the generation of a most appropriate problem-dependent topology for a given space in the more general context of chaos in Sec. 4.2.

In concluding this introduction to generalized inverses we note that the inverse  $G$  of  $f$  comes very close to being a right inverse: thus even though  $AG \neq \mathbf{1}_2$  its row-reduced form

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

is to be compared with the corresponding less satisfactory

$$\begin{pmatrix} 1 & 0 & 2 & -1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

representation of  $AG_{\text{MP}}$ .

### 3. Multifunctional Extension of Function Spaces

The previous section has considered the solution of ill-posed problems as multifunctions and has shown how this solution may be constructed. Here we introduce the multifunction space  $\text{Multi}_1(X)$  as the first step toward obtaining a smallest dense extension  $\text{Multi}(X)$  of the function space  $\text{Map}(X)$ .  $\text{Multi}_1(X)$  is basic to our theory of chaos [Sengupta & Ray, 2000] in the sense that a chaotic state of a system can be fully described by such an indeterminate multifunctional state. In fact, multifunctions also enter in a natural way in describing the spectrum of nonlinear functions that we consider in Sec. 6; this is required to complete the construction of the smallest extension  $\text{Multi}(X)$  of the function space  $\text{Map}(X)$ . The main tool in obtaining the

space  $\text{Multi}_|(X)$  from  $\text{Map}(X)$  is a generalization of the technique of pointwise convergence of continuous functions to (discontinuous) functions. In the analysis below, we consider nets instead of sequences as the spaces concerned, like the topology of pointwise convergence, may not be first countable, Appendix A.1.

### 3.1. Graphical convergence of a net of functions

Let  $(X, \mathcal{U})$  and  $(Y, \mathcal{V})$  be Hausdorff spaces and  $(f_\alpha)_{\alpha \in \mathbb{D}} : X \rightarrow Y$  be a net of piecewise continuous functions, not necessarily with the same domain or range, and suppose that for each  $\alpha \in \mathbb{D}$  there is a finite set  $I_\alpha = \{1, 2, \dots, P_\alpha\}$  such that  $f_\alpha^-$  has  $P_\alpha$  functional branches possibly with different domains; obviously  $I_\alpha$  is a singleton iff  $f$  is a injective. For each  $\alpha \in \mathbb{D}$ , define functions  $(g_{\alpha i})_{i \in I_\alpha} : Y \rightarrow X$  such that

$$f_\alpha g_{\alpha i} f_\alpha = f_{\alpha i}^I \quad i = 1, 2, \dots, P_\alpha,$$

where  $f_{\alpha i}^I$  is a basic injective branch of  $f_\alpha$  on some subset of its domain:  $g_{\alpha i} f_{\alpha i}^I = 1_X$  on  $\mathcal{D}(f_{\alpha i}^I)$ ,  $f_{\alpha i}^I g_{\alpha i} = 1_Y$  on  $\mathcal{D}(g_{\alpha i})$  for each  $i \in I_\alpha$ . The use of nets and filters is dictated by the fact that we do not assume  $X$  and  $Y$  to be first countable. In the application to the theory of dynamical systems that follows,  $X$  and  $Y$  are compact subsets of  $\mathbb{R}$  when the use of sequences suffice.

In terms of the residual and cofinal subsets  $\text{Res}(\mathbb{D})$  and  $\text{Cof}(\mathbb{D})$  of a directed set  $\mathbb{D}$  (Definition A.1.7), with  $x$  and  $y$  in the equations below being taken to belong to the required domains, define subsets  $\mathcal{D}_-$  of  $X$  and  $\mathcal{R}_-$  of  $Y$  as

$$\mathcal{D}_- = \{x \in X : ((f_\nu(x))_{\nu \in \mathbb{D}} \text{ converges in } (Y, \mathcal{V}))\} \tag{30}$$

$$\mathcal{R}_- = \{y \in Y : (\exists i \in I_\nu)((g_{\nu i}(y))_{\nu \in \mathbb{D}} \text{ converges in } (X, \mathcal{U}))\} \tag{31}$$

Thus,

$\mathcal{D}_-$  is the set of points of  $X$  on which the values of a given net of functions  $(f_\alpha)_{\alpha \in \mathbb{D}}$  converge pointwise in  $Y$ . Explicitly, this is the subset of  $X$  on which subnets<sup>19</sup> in  $\text{Map}(X, Y)$  combine to form a net of functions that converge pointwise to a limit function  $F : \mathcal{D}_- \rightarrow Y$ .

$\mathcal{R}_-$  is the set of points of  $Y$  on which the values of the nets in  $X$  generated by the injective branches

of  $(f_\alpha)_{\alpha \in \mathbb{D}}$  converge pointwise in  $Y$ . Explicitly, this is the subset of  $Y$  on which subnets of injective branches of  $(f_\alpha)_{\alpha \in \mathbb{D}}$  in  $\text{Map}(Y, X)$  combine to form a net of functions that converge pointwise to a family of limit functions  $G : \mathcal{R}_- \rightarrow X$ . Depending on the nature of  $(f_\alpha)_{\alpha \in \mathbb{D}}$ , there may be more than one  $\mathcal{R}_-$  with a corresponding family of limit functions on each of them. To simplify the notation, we will usually let  $G : \mathcal{R}_- \rightarrow X$  denote all the limit functions on all the sets  $\mathcal{R}_-$ .

If we consider cofinal rather than residual subsets of  $\mathbb{D}$  then corresponding  $\mathcal{D}_+$  and  $\mathcal{R}_+$  can be expressed as

$$\mathcal{D}_+ = \{x \in X : ((f_\nu(x))_{\nu \in \text{Cof}(\mathbb{D})} \text{ converges in } (Y, \mathcal{V}))\} \tag{32}$$

$$\mathcal{R}_+ = \{y \in Y : (\exists i \in I_\nu)((g_{\nu i}(y))_{\nu \in \text{Cof}(\mathbb{D})} \text{ converges in } (X, \mathcal{U}))\}. \tag{33}$$

It is to be noted that the conditions  $\mathcal{D}_+ = \mathcal{D}_-$  and  $\mathcal{R}_+ = \mathcal{R}_-$  are necessary and sufficient for the Kuratowski convergence to exist. Since  $\mathcal{D}_+$  and  $\mathcal{R}_+$  differ from  $\mathcal{D}_-$  and  $\mathcal{R}_-$  only in having cofinal subsets of  $D$  replaced by residual ones, and since residual sets are also cofinal, it follows that  $\mathcal{D}_- \subseteq \mathcal{D}_+$  and  $\mathcal{R}_- \subseteq \mathcal{R}_+$ . The sets  $\mathcal{D}_-$  and  $\mathcal{R}_-$  serve for the convergence of a net of functions just as  $\mathcal{D}_+$  and  $\mathcal{R}_+$  are for the convergence of subnets of the nets (*adherence*). The latter sets are needed when subsequences are to be considered as sequences in their own right as, for example, in dynamical systems theory in the case of  $\omega$ -limit sets.

As an illustration of these definitions, consider the sequence of injective functions on the interval  $[0, 1]$   $f_n(x) = 2^n x$ , for  $x \in [0, 1/2^n]$ ,  $n = 0, 1, 2, \dots$ . Then  $\mathbb{D}_{0,2}$  is the set  $\{0, 1, 2\}$  and only  $\mathbb{D}_0$  is eventual in  $\mathbb{D}$ . Hence  $\mathcal{D}_-$  is the single point set  $\{0\}$ . On the other hand  $\mathbb{D}_y$  is eventual in  $\mathbb{D}$  for all  $y$  and  $\mathcal{R}_-$  is  $[0, 1]$ .

**Definition 3.1** (Graphical Convergence of a net of functions). A net of functions  $(f_\alpha)_{\alpha \in \mathbb{D}} : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  is said to converge graphically if either  $\mathcal{D}_- \neq \emptyset$  or  $\mathcal{R}_- \neq \emptyset$ ; in this case let  $F : \mathcal{D}_- \rightarrow Y$  and  $G : \mathcal{R}_- \rightarrow X$  be the entire collection of limit functions. Because of the assumed Hausdorffness of  $X$  and  $Y$ , these limits are well defined.

The graph of the graphical limit  $\mathcal{M}$  of the net  $(f_\alpha) : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  denoted by  $f_\alpha \xrightarrow{\mathcal{G}} \mathcal{M}$ , is the

<sup>19</sup>A subnet is the generalized uncountable equivalent of a subsequence; for the technical definition, see Appendix A.1.



subset of  $\mathcal{D}_- \times \mathcal{R}_-$  that is the union of the graphs of the function  $F$  and the multifunction  $G^-$

$$\mathbf{G}_{\mathcal{M}} = \mathbf{G}_F \cup \mathbf{G}_{G^-}$$

where

$$\mathbf{G}_{G^-} = \{(x, y) \in X \times Y : (y, x) \in \mathbf{G}_G \subseteq Y \times X\}.$$

**Begin Tutorial 6: Graphical Convergence**

The following two examples are basic to the understanding of the graphical convergence of functions to multifunctions and were the examples that motivated our search of an acceptable technique that did not require vertical portions of limit relations to disappear simply because they were non-functions: the disturbing question that needed an answer was how not to mathematically sacrifice these extremely significant physical components of the limiting correspondences. Furthermore, it appears to be quite plausible to expect a physical interaction between two

spaces  $X$  and  $Y$  to be a consequence of both the direct interaction represented by  $f : X \rightarrow Y$  and also the inverse interaction  $f^- : Y \dashrightarrow X$ , and our formulation of pointwise biconvergence is a formalization of this idea. Thus the basic examples (1) and (2) below produce multifunctions instead of discontinuous functions that would be obtained by the usual pointwise limit.

**Example 3.1**

$$(1) \quad f_n(x) = \begin{cases} 0, & -1 \leq x \leq 0 \\ nx, & 0 \leq x \leq \frac{1}{n} \\ 1, & \frac{1}{n} \leq x \leq 1 \end{cases} : [-1, 1] \rightarrow [0, 1]$$

$$g_n(y) = \frac{y}{n} : [0, 1] \rightarrow \left[0, \frac{1}{n}\right]$$

Then

$$F(x) = \begin{cases} 0, & -1 \leq x \leq 0 \\ 1, & 0 < x \leq 1 \end{cases} \quad \text{on } \mathcal{D}_- = \mathcal{D}_+ = [-1, 0] \cup (0, 1]$$

$$G(y) = 0 \quad \text{on } \mathcal{R}_- = [0, 1] = \mathcal{R}_+.$$

The graphical limit is  $([-1, 0], 0) \cup (0, [0, 1]) \cup ((0, 1], 1)$ .

(2)  $f_n(x) = nx$  for  $x \in [0, 1/n]$  gives  $g_n(y) = y/n : [0, 1] \rightarrow [0, 1/n]$ . Then

$$F(x) = 0 \quad \text{on } \mathcal{D}_- = \{0\} = \mathcal{D}_+,$$

$$G(y) = 0 \quad \text{on } \mathcal{R}_- = [0, 1] = \mathcal{R}_+.$$

The graphical limit is  $(0, [0, 1])$ .

In these examples that we consider to be the prototypes of graphical convergence of functions to multifunctions,  $G(y) = 0$  on  $\mathcal{R}_-$  because  $g_n(y) \rightarrow 0$  for all  $y \in \mathcal{R}_-$ . Compare the graphical multifunctional limits with the corresponding usual pointwise functional limits characterized by discontinuity at  $x = 0$ . Two more examples from Sengupta and Ray [2000] that illustrate this new convergence principle tailored specifically to capture one-to-many relations are shown in Fig. 7 which also provides an example in Fig. 7(c) of a function whose iterates do

not converge graphically because in this case both the sets  $\mathcal{D}_-$  and  $\mathcal{R}_-$  are empty. The power of graphical convergence in capturing multifunctional limits is further demonstrated by the example of the sequence  $(\sin n\pi x)_{n=1}^\infty$  that converges to 0 both 1-integrally and test-functionally, Eqs. (3) and (4).

It is necessary to understand how the concepts of *eventually in* and *frequently in* of Appendix A.2 apply in examples (a) and (b) of Fig. 7. In these two examples we have two subsequences one each for the even indices and the other for the odd. For a point-to-point functional relation, this would mean that the sequence frequents the adherence set  $\text{adh}(x)$  of the sequence  $(x_n)$  but does not converge anywhere as it is not eventually in every neighborhood of any point. For a multifunctional limit however it is possible, as demonstrated by these examples, for the subsequences to be eventually in every neighborhood of certain *subsets* common to the eventual limiting sets of the subsequences; this intersection of the subsequential limits is now

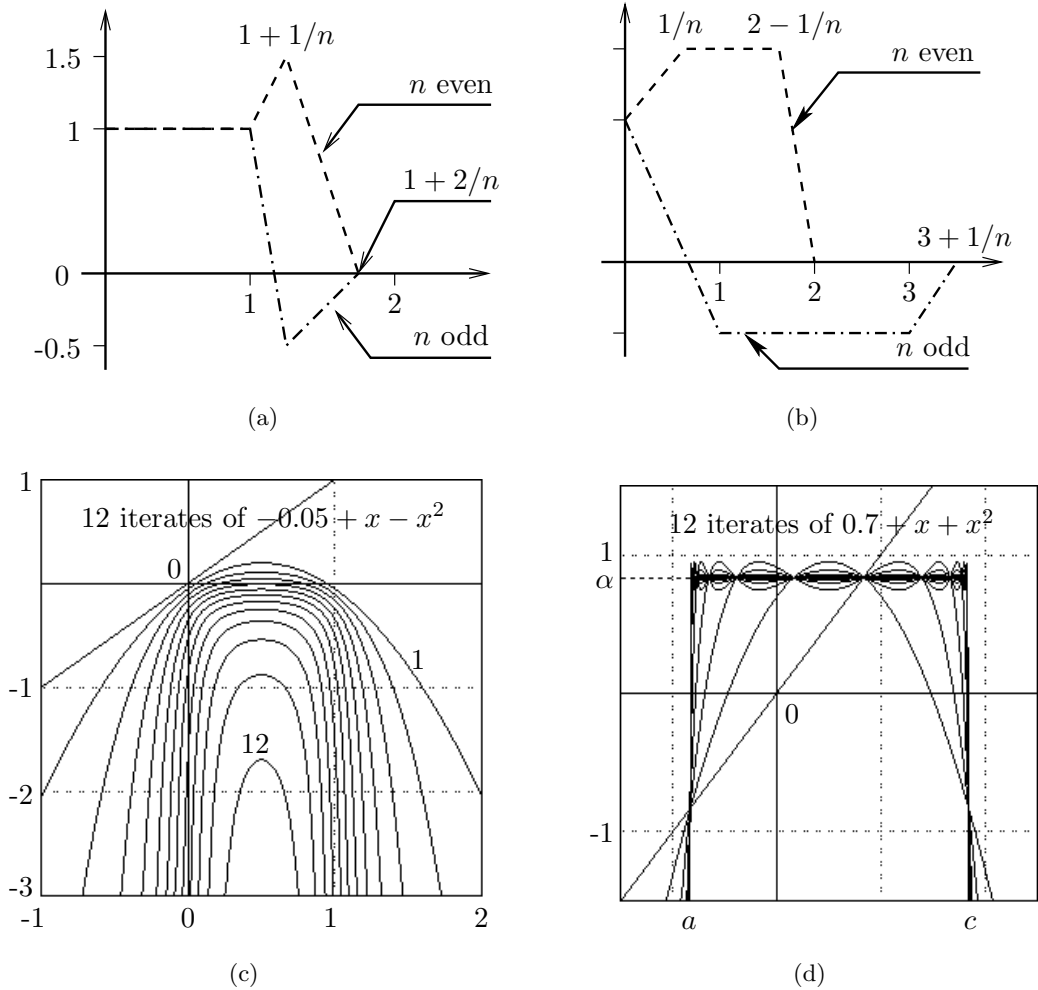


Fig. 7. The graphical limits are: (a)  $F(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1 \\ 0 & \text{for } 1 < x \leq 2 \end{cases}$  on  $\mathcal{D}_- = [0, 1] \cup (1, 2]$ , and  $G(y) = 1$  on  $\mathcal{R}_- = [0, 1]$ . Also

$$G = \begin{cases} 1 & \text{on } \mathcal{R}_+ = [0, 3/2] \\ 1 & \text{on } \mathcal{R}_+ = [-1/2, 1] \end{cases}$$

(b)  $F(x) = 1$  on  $\mathcal{D}_- = \{0\}$  and  $G(y) = 0$  on  $\mathcal{R}_- = \{1\}$ . Also  $F(x) = -1/2, 0, 1, 3/2$  respectively on  $\mathcal{D}_+ = (0, 3], \{2\}, \{0\}, (0, 2)$  and  $G(y) = 0, 0, 2, 3$  respectively on  $\mathcal{R}_+ = (-1/2, 1], [1, 3/2), [0, 3/2), [-1/2, 0)$ .

(c) For  $f(x) = -0.05 + x - x^2$ , no graphical limit as  $\mathcal{D}_- = \emptyset = \mathcal{R}_-$ .

(d) For  $f(x) = 0.7 + x - x^2$ ,  $F(x) = \alpha$  on  $\mathcal{D}_- = [a, c]$ ,  $G_1(y) = a$  and  $G_2(y) = c$  on  $\mathcal{R}_- = (-\infty, \alpha]$ . Notice how the two fixed points and their equivalent images define the converged limit rectangular multi. As in example (1) one has  $\mathcal{D}_- = \mathcal{D}_+$ ; also  $\mathcal{R}_- = \mathcal{R}_+$ .

defined to be the limit of the original sequence. A similar situation obtains, for example, in the solution of simultaneous equations: The solution of the equation  $a_{11}x_1 + a_{12}x_2 = b_1$  for one of the variables  $x_2$  say with  $a_{12} \neq 0$ , is the set represented by the straight line  $x_2 = m_1x_1 + c_1$  for all  $x_1$  in its domain, while for a different set of constants  $a_{21}, a_{22}$  and  $b_2$  the solution is the entirely different set  $x_2 = m_2x_1 + c_2$ , under the assumption that  $m_1 \neq m_2$  and  $c_1 \neq c_2$ . Thus even though the individual equations (subsequences) of the simultane-

ous set of equations (sequence) may have distinct solutions (limits), the solution of the equations is their common point of intersection.

Considered as sets in  $X \times Y$ , the discussion of convergence of a sequence of graphs  $f_n : X \rightarrow Y$  would be incomplete without a mention of the convergence of a sequence of sets under the Hausdorff metric that is so basic in the study of fractals. In this case, one talks about the convergence of a sequence of compact subsets of the metric space  $\mathbb{R}^n$  so that the sequences, as also the limit points that

are the fractals, are compact subsets of  $\mathbb{R}^n$ . Let  $\mathcal{K}$  denote the collection of all nonempty compact subsets of  $\mathbb{R}^n$ . Then the Hausdorff metric  $d_H$  between two sets on  $\mathcal{K}$  is defined to be

$$d_H(E, F) = \max\{\delta(E, F), \delta(F, E)\} \quad E, F \in \mathcal{K},$$

where

$$\delta(E, F) = \max_{x \in E} \min_{y \in F} \|x - y\|_2$$

is  $\delta(E, F)$  is the non-symmetric 2-norm in  $\mathbb{R}^n$ . The power and utility of the Hausdorff distance is best understood in terms of the dilations  $E + \varepsilon := \bigcup_{x \in E} D_\varepsilon(x)$  of a subset  $E$  of  $\mathbb{R}^n$  by  $\varepsilon$  where  $D_\varepsilon(x)$  is a closed ball of radius  $\varepsilon$  at  $x$ ; physically a dilation of  $E$  by  $\varepsilon$  is a closed  $\varepsilon$ -neighborhood of  $E$ . Then a fundamental property of  $d_H$  is that  $d_H(E, F) \leq \varepsilon$  iff both  $E \subseteq F + \varepsilon$  and  $F \subseteq E + \varepsilon$  hold simultaneously which leads [Falconer, 1990] to the interesting consequence that

If  $(F_n)_{n=1}^\infty$  and  $F$  are nonempty compact sets, then  $\lim_{n \rightarrow \infty} F_n = F$  in the Hausdorff metric iff  $F_n \subseteq F + \varepsilon$  and  $F \subseteq F_n + \varepsilon$  eventually. Furthermore if  $(F_n)_{n=1}^\infty$  is a decreasing sequence of elements of a filter-base in  $\mathbb{R}^n$ , then the nonempty and compact limit set  $F$  is given by

$$\lim_{n \rightarrow \infty} F_n = F = \bigcap_{n=1}^\infty F_n.$$

Note that since  $\mathbb{R}^n$  is Hausdorff, the assumed compactness of  $F_n$  ensures that they are also closed in  $\mathbb{R}^n$ ;  $F$ , therefore, is just the adherent set of the filter-base. In the deterministic algorithm for the generation of fractals by the so-called iterated function system (IFS) approach,  $F_n$  is the inverse image by the  $n$ th iterate of a non-injective function  $f$  having a finite number of injective branches and converging graphically to a multifunction. Under the conditions stated above, the Hausdorff metric ensures convergence of any class of compact subsets in  $\mathbb{R}^n$ . It appears eminently plausible that our multifunctional graphical convergence on  $\text{Map}(\mathbb{R}^n)$  implies Hausdorff convergence on  $\mathbb{R}^n$ : in fact pointwise biconvergence involves simultaneous convergence of image and preimage nets on  $Y$  and  $X$ , respectively. Thus confining ourselves to the simpler case of pointwise convergence, if  $(f_\alpha)_{\alpha \in \mathbb{D}}$  is a net of functions in  $\text{Map}(X, Y)$ , then the following theorem expresses the link between convergence in  $\text{Map}(X, Y)$  and in  $Y$ .

**Theorem 3.1.** *A net of functions  $(f_\alpha)_{\alpha \in \mathbb{D}}$  converges to a function  $f$  in  $(\text{Map}(X, Y), \mathcal{T})$  in the*

*topology of pointwise convergence iff  $(f_\alpha)$  converges pointwise to  $f$  in the sense that  $f_\alpha(x) \rightarrow f(x)$  in  $Y$  for every  $x$  in  $X$ .*

*Proof. Necessity.* First consider  $f_\alpha \rightarrow f$  in  $(\text{Map}(X, Y), \mathcal{T})$ . For an open neighborhood  $V$  of  $f(x)$  in  $Y$  with  $x \in X$ , let  $B(x; V)$  be a local neighborhood of  $f$  in  $(\text{Map}(X, Y), \mathcal{T})$ , see Eq. (A.6) in Appendix A.1. By assumption of convergence,  $(f_\alpha)$  must eventually be in  $B(x; V)$  implying that  $f_\alpha(x)$  is eventually in  $V$ . Hence  $f_\alpha(x) \rightarrow f(x)$  in  $Y$ .

*Sufficiency.* Conversely, if  $f_\alpha(x) \rightarrow f(x)$  in  $Y$  for every  $x \in X$ , then for a finite collection of points  $(x_i)_{i=1}^I$  of  $X$  ( $X$  may itself be uncountable) and corresponding open sets  $(V_i)_{i=1}^I$  in  $Y$  with  $f(x_i) \in V_i$ , let  $B((x_i)_{i=1}^I; (V_i)_{i=1}^I)$  be an open neighborhood of  $f$ . From the assumed pointwise convergence  $f_\alpha(x_i) \rightarrow f(x_i)$  in  $Y$  for  $i = 1, 2, \dots, I$ , it follows that  $(f_\alpha(x_i))$  is eventually in  $V_i$  for every  $(x_i)_{i=1}^I$ . Because  $\mathbb{D}$  is a directed set, the existence of a residual applicable globally for all  $i = 1, 2, \dots, I$  is assured leading to the conclusion that  $f_\alpha(x_i) \in V_i$  eventually for every  $i = 1, 2, \dots, I$ . Hence  $f_\alpha \in B((x_i)_{i=1}^I; (V_i)_{i=1}^I)$  eventually; this completes the demonstration that  $f_\alpha \rightarrow f$  in  $(\text{Map}(X, Y), \mathcal{T})$ , and thus of the proof. ■

### End Tutorial 6

---

### 3.2. The extension $\text{Multi}_|(X, Y)$ of $\text{Map}(X, Y)$

In this section we show how the topological treatment of pointwise convergence of functions to functions given in Example A.1.1 of Appendix 1 can be generalized to generate the boundary  $\text{Multi}_|(X, Y)$  between  $\text{Map}(X, Y)$  and  $\text{Multi}(X, Y)$ ; here  $X$  and  $Y$  are Hausdorff spaces and  $\text{Map}(X, Y)$  and  $\text{Multi}(X, Y)$  are respectively the sets of all functional and non-functional relations between  $X$  and  $Y$ . The generalization we seek defines neighborhoods of  $f \in \text{Map}(X, Y)$  to consist of those functional relations in  $\text{Multi}(X, Y)$  whose images at any point  $x \in X$  lies not only arbitrarily close to  $f(x)$  (this generates the usual topology of pointwise convergence  $\mathcal{T}_Y$  of Example A.1.1) but whose inverse images at  $y = f(x) \in Y$  contain points arbitrarily close to  $x$ . Thus the graph of  $f$  must not only lie close enough to  $f(x)$  at  $x$  in  $V$ , but must additionally be such that  $f^{-1}(y)$  has at least branch in  $U$  about  $x$ ; thus  $f$  is constrained to cling to  $f$  as the number of points on the graph of  $f$  increases

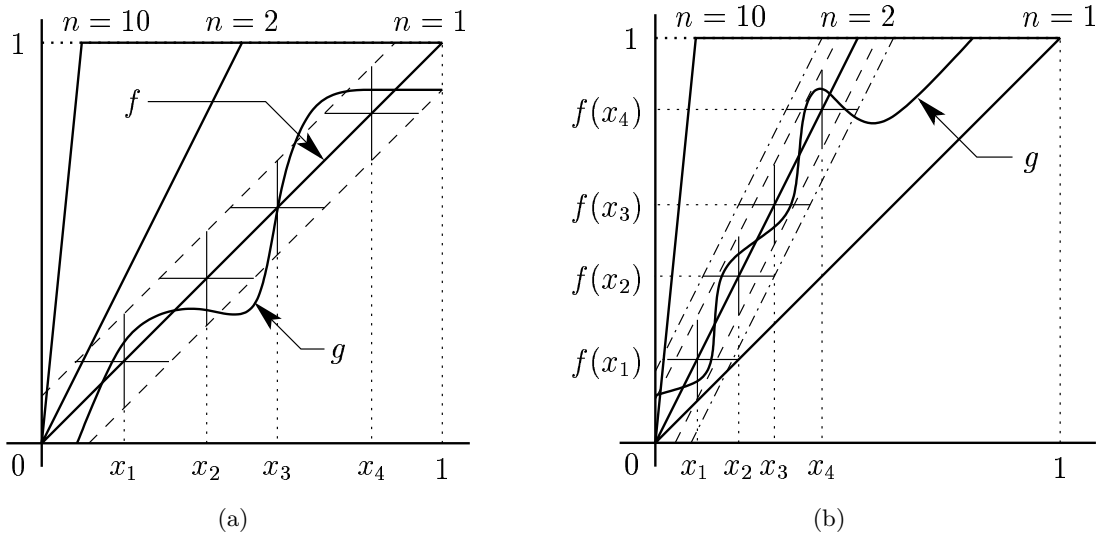


Fig. 8. The power of graphical convergence, illustrated for Example 3.1 (1), shows a local neighborhood of the functions  $x$  and  $2x$  in (a) and (b) at the four points  $(x_i)_{i=1}^4$  with corresponding neighborhoods  $(U_i)_{i=1}^4$  and  $(V_i)_{i=1}^4$  at  $(x_i, f(x_i))$  in  $\mathbb{R}$  in the  $X$  and  $Y$  directions respectively, see Eqs. (34) and (A.6) for the notations. (a) shows a function  $g$  in a pointwise neighborhood of  $f$  determined by the open sets  $V_i$ , while (b) shows  $g$  in a graphical neighborhood of  $f$  due to both  $U_i$  and  $V_i$ . A comparison of these figures demonstrates how the graphical neighborhood forces functions close to  $f$  remain closer to it than if they were in its pointwise neighborhood. This property is clearly visible in (a) where  $g$ , if it were to be in a graphical neighborhood of  $f$ , would be more faithful to it by having to be also in  $U_2$  and  $U_4$ . Thus in this case not only must the images  $f(x_{ij}) \xrightarrow{j} f(x_i)$  as  $V_i$  decreases, but also the preimages  $x_{ij} \xrightarrow{j} x_i$  with shrinking  $U_i$ . It is this simultaneous convergence of both images and preimages at every  $x$  that makes graphical convergence a natural candidate for multifunctional convergence of functions.

with convergence and, unlike in the situation of simple pointwise convergence, no gaps in the graph of the limit object is permitted not only, as in Example A.1.1 on the domain of  $f$ , but simultaneously on its range too. We call the resulting generated topology the *topology of pointwise biconvergence* on  $\text{Map}(X, Y)$ , to be denoted by  $\mathcal{T}$ . Thus for any given integer  $I \geq 1$ , the generalization of Eq. (A.6) gives for  $i = 1, 2, \dots, I$ , the open sets of  $(\text{Map}(X, Y), \mathcal{T})$  to be

$$\begin{aligned}
 & B((x_i), (V_i); (y_i), (U_i)) \\
 &= \{g \in \text{Map}(X, Y) : (g(x_i) \in V_i) \\
 &\wedge (g^{-1}(y_i) \cap U_i \neq \emptyset), i = 1, 2, \dots, I\}, \quad (34)
 \end{aligned}$$

where  $(x_i)_{i=1}^I, (V_i)_{i=1}^I$  are as in that example,  $(y_i)_{i=1}^I \in Y$ , and the corresponding open sets  $(U_i)_{i=1}^I$  in  $X$  are chosen arbitrarily.<sup>20</sup> A local base at  $f$ , for  $(x_i, y_i) \in \mathbf{G}_f$ , is the set of functions of (34) with  $y_i = f(x_i)$  and the collection of all local bases

$$B_\alpha = B((x_i)_{i=1}^{I_\alpha}, (V_i)_{i=1}^{I_\alpha}; (y_i)_{i=1}^{I_\alpha}, (U_i)_{i=1}^{I_\alpha}), \quad (35)$$

for every choice of  $\alpha \in \mathbb{D}$ , is a base  ${}_{\mathbb{T}}\mathcal{B}$  of  $(\text{Map}(X, Y), \mathcal{T})$ . Here the directed set  $\mathbb{D}$  is used as an indexing tool because, as pointed out in Example A.1.1, the topology of pointwise convergence is not first countable.

In a manner similar to Eq. (34), the open sets of  $(\text{Multi}(X, Y), \hat{\mathcal{T}})$ , where  $\text{Multi}(X, Y)$  are multifunctions with only countably many values in  $Y$  for every point of  $X$  (so that we exclude continuous regions from our discussion except for the “vertical lines” of  $\text{Multi}(X, Y)$ ), can be defined as

$$\begin{aligned}
 & \hat{B}((x_i), (V_i); (y_i), (U_i)) \\
 &= \{\mathcal{G} \in \text{Multi}(X, Y) : (\mathcal{G}(x_i) \cap V_i \neq \emptyset) \\
 &\wedge (\mathcal{G}^{-1}(y_i) \cap U_i \neq \emptyset)\}, \quad (36)
 \end{aligned}$$

where

$$\mathcal{G}^{-1}(y) = \{x \in X : y \in \mathcal{G}(x)\}.$$

and  $(x_i)_{i=1}^I \in \mathcal{D}(\mathcal{M}), (V_i)_{i=1}^I; (y_i)_{i=1}^I \in \mathcal{R}(\mathcal{M}), (U_i)_{i=1}^I$  are chosen as in the above. The topology  $\hat{\mathcal{T}}$  of  $\text{Multi}(X, Y)$  is generated by the collection of

<sup>20</sup>Equation (34) is essentially the intersection of the pointwise topologies (A.6) due to  $f$  and  $f^{-1}$ .

all local bases  $\hat{B}_\alpha$  for every choice of  $\alpha \in \mathbb{D}$ , and it is not difficult to see from Eqs. (34) and (36), that the restriction  $\hat{T}|_{\text{Map}(X, Y)}$  of  $\hat{T}$  to  $\text{Map}(X, Y)$  is just  $\mathcal{T}$ .

Henceforth  $\hat{T}$  and  $\mathcal{T}$  will be denoted by the same symbol  $\mathcal{T}$ , and convergence in the topology of pointwise biconvergence in  $(\text{Multi}(X, Y), \mathcal{T})$  will be denoted by  $\rightrightarrows$ , with the notation being derived from Theorem 3.1.

**Definition 3.2** (Functionization of a multifunction). A net of functions  $(f_\alpha)_{\alpha \in \mathbb{D}}$  in  $\text{Map}(X, Y)$  converges in  $(\text{Multi}(X, Y), \mathcal{T})$ ,  $f_\alpha \rightrightarrows \mathcal{M}$ , if it biconverges pointwise in  $(\text{Map}(X, Y), \mathcal{T}^*)$ . Such a net of functions will be said to be a functionization of  $\mathcal{M}$ .

**Theorem 3.2.** Let  $(f_\alpha)_{\alpha \in \mathbb{D}}$  be a net of functions in  $\text{Map}(X, Y)$ . Then

$$f_\alpha \xrightarrow{\mathbf{G}} \mathcal{M} \Leftrightarrow f_\alpha \rightrightarrows \mathcal{M}.$$

*Proof.* If  $(f_\alpha)$  converges graphically to  $\mathcal{M}$  then either  $\mathcal{D}_-$  or  $\mathcal{R}_-$  is non-empty; let us assume both of them to be so. Then the sequence of functions  $(f_\alpha)$  converges pointwise to a function  $F$  on  $\mathcal{D}_-$  and to functions  $G$  on  $\mathcal{R}_-$ , and the local basic neighborhoods of  $F$  and  $G$  generate the topology of pointwise biconvergence.

Conversely, for pointwise biconvergence on  $X$  and  $Y$ ,  $\mathcal{R}_-$  and  $\mathcal{D}_-$  must be non-empty. ■

Observe that the boundary of  $\text{Map}(X, Y)$  in the topology of pointwise biconvergence is a “line parallel to the  $Y$ -axis”. We denote this closure of  $\text{Map}(X, Y)$  as

**Definition 3.3.**  $\text{Multi}_|((X, Y), \mathcal{T}) = \text{Cl}(\text{Map}((X, Y), \mathcal{T}))$ .

The sense in which  $\text{Multi}_|(X, Y)$  is the smallest closed topological extension of  $M = \text{Map}(X, Y)$  is the following, refer to Theorems A.1.4 and its proof. Let  $(M, \mathcal{T}_0)$  be a topological space and suppose that

$$\hat{M} = M \cup \{\hat{m}\}$$

is obtained by adjoining an extra point to  $M$ ; here  $M = \text{Map}(X, Y)$  and  $\hat{m} \in \text{Cl}(M)$  is the multifunctional limit in  $\hat{M} = \text{Multi}_|(X, Y)$ . Treat all open sets of  $M$  generated by local bases of the type (35) with finite intersection property as a filter-base  ${}_F\mathcal{B}$  on  $X$  that induces a filter  $\mathcal{F}$  on  $M$  (by forming su-

persets of all elements of  ${}_F\mathcal{B}$ ; see Appendix A.1) and thereby the filter-base

$${}_F\hat{\mathcal{B}} = \{\hat{B} = B \cup \{\hat{m}\} : B \in {}_F\mathcal{B}\}$$

on  $\hat{M}$ ; this filter-base at  $m$  can also be obtained independently from Eq. (36). Obviously  ${}_F\hat{\mathcal{B}}$  is an extension of  ${}_F\mathcal{B}$  on  $\hat{M}$  and  ${}_F\mathcal{B}$  is the filter induced on  $M$  by  ${}_F\hat{\mathcal{B}}$ . We may also consider the filter-base to be a topological base on  $M$  that defines a coarser topology  $\mathcal{T}$  on  $M$  (through all unions of members of  ${}_F\mathcal{B}$ ) and hence the topology

$$\hat{\mathcal{T}} = \{\hat{G} = G \cup \{\hat{m}\} : G \in \mathcal{T}\}$$

on  $\hat{M}$  to be the topology associated with  $\hat{\mathcal{F}}$ . A finer topology on  $\hat{M}$  may be obtained by adding to  $\hat{\mathcal{T}}$  all the discarded elements of  $\mathcal{T}_0$  that do not satisfy FIP. It is clear that  $\hat{m}$  is on the boundary of  $M$  because every neighborhood of  $\hat{m}$  intersects  $M$  by construction; thus  $(M, \mathcal{T})$  is dense in  $(\hat{M}, \hat{\mathcal{T}})$  which is the required topological extension of  $(M, \mathcal{T})$ .

In the present case, a filter-base at  $f \in \text{Map}(X, Y)$  is the neighborhood system  ${}_F\mathcal{B}_f$  at  $f$  given by decreasing sequences of neighborhoods  $(V_k)$  and  $(U_k)$  of  $f(x)$  and  $x$ , respectively, and the filter  $\hat{\mathcal{F}}$  is the neighborhood filter  $\mathcal{N}_f \cup G$  where  $G \in \text{Multi}_|(X, Y)$ . We shall present an alternate, and perhaps more intuitively appealing, description of graphical convergence based on the adherence set of a filter in Sec. 4.1.

As more serious examples of the graphical convergence of a net of functions to multifunction than those considered above, Fig. 9 shows the first four iterates of the tent map

$$t(x) = \begin{cases} 2x, & 0 \leq x < \frac{1}{2} \\ 2(1-x), & \frac{1}{2} \leq x \leq 1 \end{cases} \quad (t^1 = t).$$

defined on  $[0, 1]$  and the sine map  $f_n = |\sin(2^{n-1}\pi x)|$ ,  $n = 1, \dots, 4$  with domain  $[0, 1]$ .

These examples illustrate the important generalization that *periodic points may be replaced by the more general equivalence classes* where a sequence of functions converges graphically; this generalization based on the ill-posed interpretation of dynamical systems is significant for non-iterative systems as in second example above. The equivalence classes of the tent map for its two fixed points 0 and 2/3 generated by the first four iterates are

$$[0]_4 = \left\{ 0, \frac{1}{8}, \frac{1}{4}, \frac{3}{8}, \frac{1}{2}, \frac{5}{8}, \frac{3}{4}, \frac{7}{8}, 1 \right\}$$

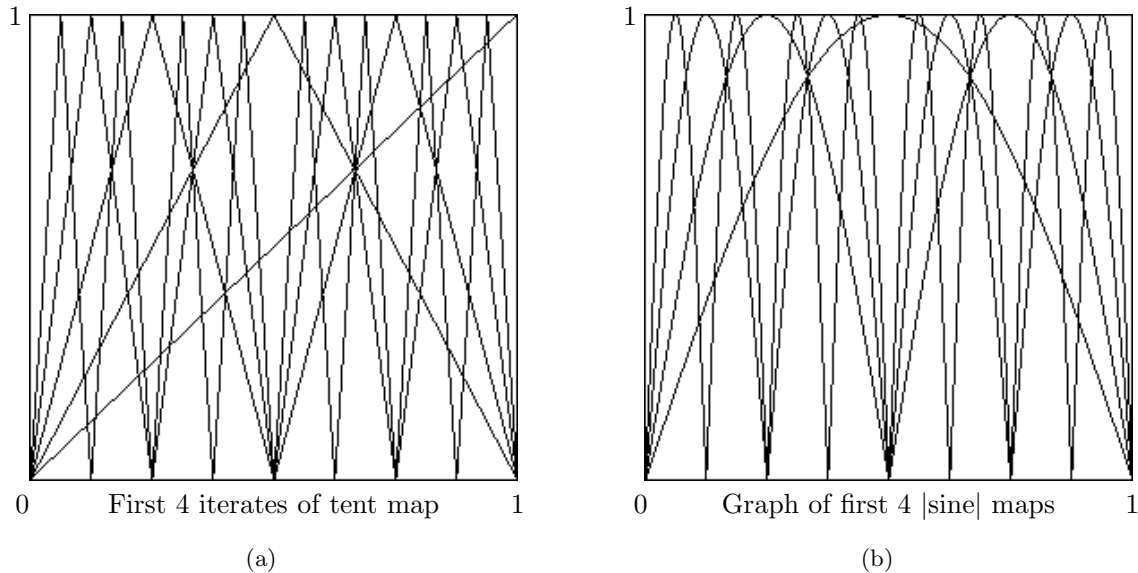


Fig. 9. The first four iterates of (a) tent and (b)  $|\sin(2^{n-1}\pi x)|$  maps show the formal similarity of the dynamics of these functions. It should be noted, as shown in Fig. 7, that although  $\sin(n\pi x)_{n=1}^\infty$  fails to converge at any point other than 0 and 1, the subsequence  $\sin(2^{n-1}\pi x)_{n=1}^\infty$  does converge graphically on a set dense in  $[0, 1]$ .

$$\left[ \frac{2}{3} \right]_4 = \left\{ c, \frac{1}{8} \mp c, \frac{1}{4} \mp c, \frac{3}{8} \mp c, \frac{1}{2} \mp c, \frac{5}{8} \mp c, \frac{3}{4} \mp c, \frac{7}{8} \mp c, 1 - c \right\}$$

where  $c = 1/24$ . If the moduli of the slopes of the graphs passing through these equivalent fixed points are greater than 1 then the graphs converge to multifunctions and when these slopes are less than 1 the corresponding graphs converge to constant

functions. It is to be noted that the number of equivalent fixed points in a class increases with the number of iterations  $k$  as  $2^{k-1} + 1$ ; this increase in the degree of ill-posedness is typical of discrete chaotic systems and can be regarded as a paradigm of chaos generated by the convergence of a family of functions.

The  $m$ th iterate  $t^m$  of the tent map has  $2^m$  fixed points corresponding to the  $2^m$  injective branches of  $t^m$

$$x_{mj} = \begin{cases} \frac{j-1}{2^m-1}, & j = 1, 3, \dots, (2^m-1) \\ \frac{j}{2^m+1}, & j = 2, 4, \dots, 2^m \end{cases} \quad t^m(x_{mj}) = x_{mj}, j = 1, 2, \dots, 2^m.$$

Let  $X_m$  be the collection of these  $2^m$  fixed points (thus  $X_1 = \{0, 2/3\}$ ), and denote by  $[X_m]$  the set of the equivalent points, one coming from each of the injective branches, for each of the fixed points: thus

$$\mathcal{D}_- = [X_1] = \left\{ [0], \left[ \frac{2}{3} \right] \right\}$$

$$[X_2] = \left\{ [0], \left[ \frac{2}{5} \right], \left[ \frac{2}{3} \right], \left[ \frac{4}{5} \right] \right\}$$

and  $\mathcal{D}_+ = \bigcap_{m=1}^\infty [X_m]$  is a non-empty countable set dense in  $X$  at each of which the graphs of the sequence  $(t^m)$  converge to a multifunction. New sets  $[X_n]$  will be formed by subsequences of the

higher iterates  $t^n$  for  $m = in$  with  $i = 1, 2, \dots$  where these subsequences remain fixed. For example, the fixed points  $2/5$  and  $4/5$  produced respectively by the second and fourth injective branches of  $t^2$ , are also fixed for the seventh and thirteenth branches of  $t^4$ . For the shift map  $2x \bmod(1)$  on  $[0, 1]$ ,  $\mathcal{D}_- = \{[0], [1]\}$  where  $[0] = \bigcap_{m=1}^\infty \{(i-1)/2^m : i = 1, 2, \dots, 2^m\}$  and  $[1] = \bigcap_{m=1}^\infty \{i/2^m : i = 1, 2, \dots, 2^m\}$ .

It is useful to compare the graphical convergence of  $(\sin(\pi nx))_{n=1}^\infty$  to  $[0, 1]$  at 0 and to 0 at 1 with the usual integral and test-functional convergences to 0; note that the point  $1/2$ , for example, belongs to  $\mathcal{D}_+$  and not to  $\mathcal{D}_- = \{0, 1\}$  because it is frequented by even  $n$  only. However for the sub-

sequence  $(f_{2^{m-1}})_{m \in \mathbb{Z}_+}$ ,  $1/2$  is in  $\mathcal{D}_-$  because if the graph of  $f_{2^{m-1}}$  passes through  $(1/2, 0)$  for some  $m$ , then so do the graphs for all higher values. Therefore  $[0] = \bigcap_{m=1}^{\infty} \{i/2^{m-1} : i = 0, 1, \dots, 2^{m-1}\}$  is the equivalence class of  $(f_{2^{m-1}})_{m=1}^{\infty}$  and this sequence converges to  $[-1, 1]$  on this set. Thus our extension  $\text{Multi}(X)$  is distinct from the distributional extension of function spaces with respect to test functions, and is able to correctly generate the pathological behavior of the limits that are so crucially vital in producing chaos.

#### 4. Discrete Chaotic Systems are Maximally Ill-posed

The above ideas apply to the development of a criterion for chaos in discrete dynamical systems that is based on the limiting behavior of the graphs of a sequence of functions  $(f_n)$  on  $X$ , rather than on the values that the sequence generates as is customary. For the development of the maximality of ill-posedness criterion of chaos, we need to refresh ourselves with the following preliminaries.

---

#### Resume Tutorial 5: Axiom of Choice and Zorn’s Lemma

Let us recall from the first part of this Tutorial that for nonempty subsets  $(A_\alpha)_{\alpha \in \mathbb{D}}$  of a nonempty set  $X$ , the Axiom of Choice ensures the existence of a set  $A$  such that  $A \cap A_\alpha$  consists of a single element for every  $\alpha$ . The choice axiom has far reaching consequences and a few equivalent statements, one of which the Zorn’s lemma that will be used immediately in the following, is the topic of this resumed Tutorial. The beauty of the Axiom, and of its equivalents, is that they assert the existence of mathematical objects that, in general, cannot be demonstrated and it is often believed that Zorn’s lemma is one of the most powerful tools that a mathematician has available to him that is “almost indispensable in many parts of modern pure mathematics” with significant applications in nearly all branches of contemporary mathematics. This “lemma” talks about maximal (as distinct from “maximum”) elements of a partially ordered set, a set in which some notion of  $x_1$  “preceding”  $x_2$  for two elements of the set has been defined.

A relation  $\preceq$  on a set  $X$  is said to be a *partial order* (or simply an *order*) if it is (compare with the

properties (ER1)–(ER3) of an equivalence relation, Tutorial 1)

(OR1) Reflexive, that is  $(\forall x \in X)(x \preceq x)$ .

(OR2) Antisymmetric:  $(\forall x, y \in X)(x \preceq y \wedge y \preceq x \Rightarrow x = y)$ .

(OR3) Transitive, that is  $(\forall x, y, z \in X)(x \preceq y \wedge y \preceq z \Rightarrow x \preceq z)$ . Any notion of order on a set  $X$  in the sense of one element of  $X$  preceding another should possess at least this property.

The relation is a *preorder*  $\preceq$  if it is only reflexive and transitive, that is if only (OR1) and (OR3) are true. If the hypothesis of (OR2) is also satisfied by a preorder, then this  $\preceq$  induces an equivalence relation  $\sim$  on  $X$  according to  $(x \preceq y) \wedge (y \preceq x) \Leftrightarrow x \sim y$  that evidently is actually a partial order iff  $x \sim y \Leftrightarrow x = y$ . For any element  $[x] \in X/\sim$  of the induced quotient space, let  $\leq$  denote the generated order in  $X/\sim$  so that

$$x \preceq y \Leftrightarrow [x] \leq [y];$$

then  $\leq$  is a partial order on  $X/\sim$ . If every two elements of  $X$  are *comparable*, in the sense that either  $x_1 \preceq x_2$  or  $x_2 \preceq x_1$  for all  $x_1, x_2 \in X$ , then  $X$  is said to be a *totally ordered set* or a *chain*. A totally ordered subset  $(C, \preceq)$  of a partially ordered set  $(X, \preceq)$  with the ordering induced from  $X$ , is known as a *chain in X* if

$$C = \{x \in X : (\forall c \in X)(c \preceq x \vee x \preceq c)\}. \quad (37)$$

The most important class of chains that we are concerned with in this work is that on the subsets  $\mathcal{P}(X)$  of a set  $(X, \subseteq)$  under the inclusion order; Eq. (37), as we shall see in what follows, defines a family of chains of nested subsets in  $\mathcal{P}(X)$ . Thus while the relation  $\preceq$  in  $\mathbb{Z}$  defined by  $n_1 \preceq n_2 \Leftrightarrow |n_1| \leq |n_2|$  with  $n_1, n_2 \in \mathbb{Z}$  preorders  $\mathbb{Z}$ , it is not a partial order because although  $-n \preceq n$  and  $n \preceq -n$  for any  $n \in \mathbb{Z}$ , it is does not follow that  $-n = n$ . A common example of partial order on a set of sets, for example on the power set  $\mathcal{P}(X)$  of a set  $X$  (see footnote 23), is the inclusion relation  $\subseteq$ : the ordered set  $\mathcal{X} = (\mathcal{P}(\{x, y, z\}), \subseteq)$  is partially ordered but not totally ordered because, for example,  $\{x, y\} \not\subseteq \{y, x\}$ , or  $\{x\}$  is not comparable to  $\{y\}$  unless  $x = y$ ; however  $C = \{\{\emptyset, \{x\}, \{x, y\}\}$  does represent one of the many possible chains of  $\mathcal{X}$ . Another useful example of partial order is the following: Let  $X$  and  $(Y, \leq)$  be sets with  $\leq$  ordering  $Y$ , and consider  $f, g \in \text{Map}(X, Y)$  with

$\mathcal{D}(f), \mathcal{D}(g) \subseteq X$ . Then

$$\begin{aligned} (\mathcal{D}(f) \subseteq \mathcal{D}(g))(f = g|_{\mathcal{D}(f)}) &\Leftrightarrow f \preceq g \\ (\mathcal{D}(f) = \mathcal{D}(g))(\mathcal{R}(f) \subseteq \mathcal{R}(g)) &\Leftrightarrow f \preceq g \\ (\forall x \in \mathcal{D}(f) = \mathcal{D}(g))(f(x) \leq g(x)) &\Leftrightarrow f \preceq g \end{aligned} \quad (38)$$

define partial orders on  $\text{Map}(X, Y)$ . In the last case, the order is not total because any two functions whose graphs cross at some point in their common domain cannot be ordered by the given relation, while in the first any  $f$  whose graph does not coincide with that of  $g$  on the common domain is not comparable to it by this relation.

Let  $(X, \preceq)$  be a partially ordered set and let  $A$  be a subset of  $X$ . An element  $a_+ \in (A, \preceq)$  is said to be a *maximal* element of  $A$  with respect to  $\preceq$  if

$$(\forall a \in (A, \preceq))(a_+ \preceq a) \Rightarrow a = a_+, \quad (39)$$

that is, iff there is no  $a \in A$  with  $a \neq a_+$  and  $a \succ a_+$ .<sup>21</sup> Expressed otherwise, this implies that an element  $a_+$  of a subset  $A \subseteq (X, \preceq)$  is maximal in  $(A, \preceq)$  iff it is true that

$$(a \preceq a_+ \in A)(\text{for every } a \in (A, \preceq) \text{ comparable to } a_+); \quad (40)$$

thus  $a_+$  in  $A$  is a maximal element of  $A$  iff it is strictly greater than every *other comparable* element of  $A$ . This of course does not mean that each element  $a$  of  $A$  satisfies  $a \preceq a_+$  because every pair of elements of a partially ordered set need not be comparable: in a totally ordered set there can be at most one maximal element. In comparison, an element  $a_\infty$  of a subset  $A \subseteq (X, \preceq)$  is *the* unique *maximum* (*largest, greatest, last*) element of  $A$  iff

$$(a \preceq a_\infty \in A)(\text{for every } a \in (A, \preceq)), \quad (41)$$

implying that  $a_\infty$  is *the* element of  $A$  that is strictly larger than every other element of  $A$ . As in the case of the maximal, although this also does not require all elements of  $A$  to be comparable to each other, it does require  $a_\infty$  to be larger than every element of  $A$ . The dual concepts of minimal and minimum can be similarly defined by essentially reversing the roles of  $a$  and  $b$  in relational expressions like  $a \preceq b$ .

The last concept needed to formalize Zorn's lemma is that of an upper bound: For a subset  $(A, \preceq)$  of a partially ordered set  $(X, \preceq)$ , an element  $u$  of  $X$  is an *upper bound of  $A$  in  $X$*  iff

$$(a \preceq u \in (X, \preceq))(\text{for every } a \in (A, \preceq)) \quad (42)$$

which requires the upper bound  $u$  to be larger than all members of  $A$ , with the corresponding lower bounds of  $A$  being defined in a similar manner. Of course, it is again not necessary that the elements of  $A$  be comparable to each other, and it should be clear from Eqs. (41) and (42) that when an upper bound of a set is in the set itself, then it is the maximum element of the set. If the upper (lower) bounds of a subset  $(A, \preceq)$  of a set  $(X, \preceq)$  has a least (greatest) element, then this smallest upper bound (largest lower bound) is called *the least upper bound* (*greatest lower bound*) or *supremum* (*infimum*) of  $A$  in  $X$ . Combining Eqs. (41) and (42) then yields

$$\begin{aligned} \sup_X A &= \{a_\leftarrow \in \Omega_A : a_\leftarrow \preceq u \forall u \in (\Omega_A, \preceq)\} \\ \inf_X A &= \{\rightarrow a \in \Lambda_A : l \preceq \rightarrow a \forall l \in (\Lambda_A, \preceq)\} \end{aligned} \quad (43)$$

where  $\Omega_A = \{u \in X : (\forall a \in A)(a \preceq u)\}$  and  $\Lambda_A = \{l \in X : (\forall a \in A)(l \preceq a)\}$  are the sets of all upper and lower bounds of  $A$  in  $X$ . Equation (43) may be expressed in the equivalent but more transparent form as

$$\begin{aligned} a_\leftarrow = \sup_X A &\Leftrightarrow (a \in A \Rightarrow a \preceq a_\leftarrow) \\ &\wedge (a_0 \prec a_\leftarrow \Rightarrow a_0 \prec b \preceq a_\leftarrow \text{ for some } b \in A) \\ \rightarrow a = \inf_X A &\Leftrightarrow (a \in A \Rightarrow \rightarrow a \preceq a) \\ &\wedge (\rightarrow a \prec a_1 \Rightarrow \rightarrow a \preceq b \prec a_1 \text{ for some } b \in A) \end{aligned} \quad (44)$$

to imply that  $a_\leftarrow$  ( $\rightarrow a$ ) is *the* upper (lower) bound of  $A$  in  $X$  which precedes (succeeds) every other upper (lower) bound of  $A$  in  $X$ . Notice that uniqueness in the definitions above is a direct consequence of the uniqueness of greatest and least elements of a set. It must be noted that whereas maximal and maximum are properties of the particular subset and have nothing to do with anything outside it, upper and lower bounds of a set are defined only with respect to a superset that may contain it.

The following example, beside being useful in Zorn's lemma, is also of great significance in fixing some of the basic ideas needed in our future arguments involving classes of sets ordered by the inclusion relation.

**Example 4.1.** Let  $\mathcal{X} = \mathcal{P}(\{a, b, c\})$  be ordered by the inclusion relation  $\subseteq$ . The subset  $\mathcal{A} = \mathcal{P}(\{a, b, c\}) - \{a, b, c\}$  has three maximals  $\{a, b\}$ ,  $\{b, c\}$  and  $\{c, a\}$  but no maximum as there is no

<sup>21</sup>If  $\preceq$  is an order relation in  $X$  then the *strict relation*  $\prec$  in  $X$  corresponding to  $\preceq$ , given by  $x \prec y \Leftrightarrow (x \preceq y) \wedge (x \neq y)$ , is *not an order relation* because unlike  $\preceq$ ,  $\prec$  is not reflexive even though it is both transitive and asymmetric.



$A_\infty \in \mathcal{A}$  satisfying  $A \preceq A_\infty$  for every  $A \in \mathcal{A}$ , while  $\mathcal{P}(\{a, b, c\}) - \emptyset$  the three minimals  $\{a\}$ ,  $\{b\}$  and  $\{c\}$  but no minimum. This shows that a subset of a partially ordered set may have many maximals (minimals) without possessing a maximum (minimum), but a subset has a maximum (minimum) iff this is its unique maximal (minimal). If  $\mathcal{A} = \{\{a, b\}, \{a, c\}\}$ , then every subset of the intersection of the elements of  $\mathcal{A}$ , namely  $\{a\}$  and  $\emptyset$ , are lower bounds of  $\mathcal{A}$ , and all supersets in  $\mathcal{X}$  of the union of its elements — which in this case is just  $\{a, b, c\}$  — are its upper bounds. Notice that while the maximal (minimal) and maximum (minimum) are elements of  $\mathcal{A}$ , upper and lower bounds need not be contained in their sets. In this class  $(\mathcal{X}, \subseteq)$  of subsets of a set  $X$ ,  $X_+$  is a maximal element of  $\mathcal{X}$  iff  $X_+$  is not contained in any other subset of  $X$ , while  $X_\infty$  is a maximum of  $\mathcal{X}$  iff  $X_\infty$  contains every other subset of  $X$ .

Let  $\mathcal{A} := \{A_\alpha \in \mathcal{X}\}_{\alpha \in \mathbb{D}}$  be a non-empty subclass of  $(\mathcal{X}, \subseteq)$ , and suppose that both  $\bigcup A_\alpha$  and  $\bigcap A_\alpha$  are elements of  $\mathcal{X}$ . Since each  $A_\alpha$  is  $\subseteq$ -less than  $\bigcup A_\alpha$ , it follows that  $\bigcup A_\alpha$  is an upper bound of  $\mathcal{A}$ ; this is also the smallest of all such bounds because if  $U$  is any other upper bound then every  $A_\alpha$  must precede  $U$  by Eq. (42) and therefore so must  $\bigcup A_\alpha$  (because the union of a class of subsets of a set is the smallest that contain each member of the class:  $A_\alpha \subseteq U \Rightarrow \bigcup A_\alpha \subseteq U$  for subsets  $(A_\alpha)$  and  $U$  of  $X$ ). Analogously, since  $\bigcap A_\alpha$  is  $\subseteq$ -less than each  $A_\alpha$  it is a lower bound of  $\mathcal{A}$ ; that it is the greatest of all the lower bounds  $L$  in  $\mathcal{X}$  follows because the intersection of a class of subsets is the largest that is contained in each of the subsets:  $L \subseteq A_\alpha \Rightarrow L \subseteq \bigcap A_\alpha$  for subsets  $L$  and  $(A_\alpha)$  of  $X$ . Hence the supremum and infimum of  $\mathcal{A}$  in  $(\mathcal{X}, \subseteq)$  given by

$$A_- = \sup_{(\mathcal{X}, \subseteq)} \mathcal{A} = \bigcup_{A \in \mathcal{A}} A \tag{45}$$

and 
$$\rightarrow A = \inf_{(\mathcal{X}, \subseteq)} \mathcal{A} = \bigcap_{A \in \mathcal{A}} A$$

are both elements of  $(\mathcal{X}, \subseteq)$ . Intuitively, an upper (respectively, lower) bound of  $\mathcal{A}$  in  $\mathcal{X}$  is any subset of  $\mathcal{X}$  that contains (respectively, is contained in) every member of  $\mathcal{A}$ .

The statement of Zorn’s lemma and its proof can now be completed in three stages as follows. For Theorem 4.1 below that constitutes the most significant technical first stage, let  $g$  be a function on  $(X, \preceq)$  that assigns to every  $x \in X$  an *immediate successor*  $y \in X$  such that

$$\mathcal{M}(x) = \{y \succ x : \nexists x_* \in X \text{ satisfying } x \prec x_* \prec y\}$$

are all the successors of  $x$  in  $X$  with no element of  $X$  lying strictly between  $x$  and  $y$ . Select a representative of  $\mathcal{M}(x)$  by a choice function  $f_C$  such that

$$g(x) = f_C(\mathcal{M}(x)) \in \mathcal{M}(x)$$

is an immediate successor of  $x$  chosen from the many possible in the set  $\mathcal{M}(x)$ . The basic idea in the proof of the first of the three-parts is to express the existence of a maximal element of a partially ordered set  $X$  in terms of the existence of a fixed point in the set, which follows as a contradiction of the assumed hypothesis that every point in  $X$  has an immediate successor. Our basic application of immediate successors in the following will be to classes  $\mathcal{X} \subseteq (\mathcal{P}(X), \subseteq)$  of subsets of a set  $X$  ordered by inclusion. In this case for any  $A \in \mathcal{X}$ , the function  $g$  can be taken to be the superset

$$g(A) = A \cup f_C(\mathcal{G}(A) - A), \tag{46}$$

where  $\mathcal{G}(A) = \{x \in X - A : A \cup \{x\} \in \mathcal{X}\}$

of  $A$ . Repeated application of  $g$  to  $A$  then generates a principal filter, and hence an associated sequence, based at  $A$ .

**Theorem 4.1.** *Let  $(X, \preceq)$  be a partially ordered set that satisfies*

(ST1) *There is a smallest element  $x_0$  of  $X$  which has no immediate predecessor in  $X$ .*

(ST2) *If  $C \subseteq X$  is a totally ordered subset in  $X$ , then  $c_* = \sup_X C$  is in  $X$ .*

*Then there exists a maximal element  $x_+$  of  $X$  which has no immediate successor in  $X$ .*

*Proof.* Let  $T \subseteq (X, \preceq)$  be a subset of  $X$ . If the conclusion of the theorem is false then the alternative

(ST3) Every element  $x \in T$  has an immediate successor  $g(x)$  in  $T$ <sup>22</sup>

---

<sup>22</sup>This makes  $T$ , and hence  $X$ , inductively defined infinite sets. It should be realized that (ST3) *does not mean* that every member of  $T$  is obtained from  $g$ , but only ensures that the immediate successor of any element of  $T$  is also in  $T$ . The infimum  $\rightarrow T$  of these towers satisfies the additional property of being totally ordered (and is therefore essentially a sequence or net) in  $(X, \preceq)$  to which (ST2) can be applied.

leads, as shown below, to a contradiction that can be resolved only by the conclusion of the theorem. A subset  $T$  of  $(X, \preceq)$  satisfying conditions (ST1)–(ST3) is sometimes known as an  $g$ -tower or an  $g$ -sequence: an obvious example of a tower is  $(X, \preceq)$  itself. If

$$\rightarrow T = \bigcap \{T \in \mathcal{T} : T \text{ is an } x_0 \text{-tower}\}$$

is the  $(\mathcal{P}(X), \subseteq)$ -infimum of the class  $\mathcal{T}$  of all sequential towers of  $(X, \preceq)$ , we show that this smallest sequential tower is in fact a *sequential totally ordered chain* in  $(X, \preceq)$  built from  $x_0$  by the  $g$ -function. Let the subset

$$C_T = \{c \in X : (\forall t \in \rightarrow T)(t \preceq c \vee c \preceq t)\} \subseteq X \quad (47)$$

of  $X$  be an  $g$ -chain in  $\rightarrow T$  in the sense that [cf. Eq. (37)] it is that subset of  $X$  each of whose elements is comparable with some element of  $\rightarrow T$ . The conditions (ST1)–(ST3) for  $C_T$  can be verified as follows to demonstrate that  $C_T$  is an  $g$ -tower.

- (1)  $x_0 \in C_T$ , because it is less than each  $x \in \rightarrow T$ .
- (2) Let  $c_{\leftarrow} = \sup_X C_T$  be the supremum of the chain  $C_T$  in  $X$  so that by (ST2),  $c_{\leftarrow} \in X$ . Let  $t \in \rightarrow T$ . If there is *some*  $c \in C_T$  such that  $t \preceq c$ , then surely  $t \preceq c_{\leftarrow}$ . Else,  $c \preceq t$  for *every*  $c \in C_T$  shows that  $c_{\leftarrow} \preceq t$  because  $c_{\leftarrow}$  is the smallest of all the upper bounds  $t$  of  $C_T$ . Therefore  $c_{\leftarrow} \in C_T$ .
- (3) In order to show that  $g(c) \in C$  whenever  $c \in C$  it needs to be verified that for all  $t \in \rightarrow T$ , either  $t \preceq c \Rightarrow t \preceq g(c)$  or  $c \preceq t \Rightarrow g(c) \preceq t$ . As the former is clearly obvious, we investigate the latter as follows; note that  $g(t) \in \rightarrow T$  by (ST3). The first step is to show that the subset

$$C_g = \{t \in \rightarrow T : (\forall c \in C_T)(t \preceq c \vee g(c) \preceq t)\} \quad (48)$$

of  $\rightarrow T$ , which is a chain in  $X$  (observe the inverse roles of  $t$  and  $c$  here as compared to that in Eq. (47)), is a tower: Let  $t_{\leftarrow}$  be the supremum of  $C_g$  and take  $c \in C$ . If there is *some*  $t \in C_g$  for which  $g(c) \preceq t$ , then clearly  $g(c) \preceq t_{\leftarrow}$ . Else,  $t \preceq x$  for *each*  $t \in C_g$  shows that  $t_{\leftarrow} \preceq c$  because  $t_{\leftarrow}$  is the smallest of all the upper bounds  $c$  of  $C_g$ . Hence  $t_{\leftarrow} \in C_g$ .

Property (ST3) for  $C_g$  follows from a small yet significant modification of the above arguments in which the immediate successors  $g(t)$  of  $t \in C_g$  formally replaces the supremum  $t_{\leftarrow}$  of  $C_g$ . Thus given a  $c \in C$ , if there is *some*  $t \in C_g$  for which

$g(c) \preceq t$  then  $g(c) \prec g(t)$ ; this combined with  $(c = t) \Rightarrow (g(c) = g(t))$  yields  $g(c) \preceq g(t)$ . On the other hand,  $t \prec c$  for *every*  $t \in C_g$  requires  $g(t) \preceq c$  as otherwise  $(t \prec c) \Rightarrow (c \prec g(t))$  would, from the resulting consequence  $t \prec c \prec g(t)$ , contradict the assumed hypothesis that  $g(t)$  is the immediate successor of  $t$ . Hence,  $C_g$  is a  $g$ -tower in  $X$ .

To complete the proof that  $g(c) \in C_T$ , and thereby the argument that  $C_T$  is a tower, we first note that as  $\rightarrow T$  is the smallest tower and  $C_g$  is built from it,  $C_g = \rightarrow T$  must in fact be  $\rightarrow T$  itself. From Eq. (48) therefore, for every  $t \in \rightarrow T$  either  $t \preceq g(c)$  or  $g(c) \preceq t$ , so that  $g(c) \in C_T$  whenever  $c \in C_T$ . This concludes the proof that  $C_T$  is actually the tower  $\rightarrow T$  in  $X$ .

From (ST2), the implication of the chain  $C_T$

$$C_T = \rightarrow T = C_g \quad (49)$$

being the minimal tower  $\rightarrow T$  is that the supremum  $t_{\leftarrow}$  of the totally ordered  $\rightarrow T$  in its own tower (as distinct from in the tower  $X$ : recall that  $\rightarrow T$  is a subset of  $X$ ) must be contained in itself, that is

$$\sup_{C_T}(C_T) = t_{\leftarrow} \in \rightarrow T \subseteq X. \quad (50)$$

This however leads to the contradiction from (ST3) that  $g(t_{\leftarrow})$  be an element of  $\rightarrow T$ , unless of course

$$g(t_{\leftarrow}) = t_{\leftarrow}, \quad (51)$$

which because of (49) may also be expressed equivalently as  $g(c_{\leftarrow}) = c_{\leftarrow} \in C_T$ . As the sequential totally ordered set  $\rightarrow T$  is a subset of  $X$ , Eq. (48) implies that  $t_{\leftarrow}$  is a maximal element of  $X$  which allows (ST3) to be replaced by the remarkable inverse criterion that

(ST3') If  $x \in X$  and  $w$  precedes  $x$ ,  $w \prec x$ , then  $w \in X$ , that is obviously false for a general tower  $T$ . In fact, it follows directly from Eq. (39) that under (ST3') *any*  $x_+ \in X$  is a maximal element of  $X$  iff it is a fixed point of  $g$  as given by Eq. (51). This proves the theorem and also demonstrates how, starting from a minimum element of a partially ordered set  $X$ , (ST3) can be used to generate inductively a totally ordered sequential subset of  $X$  leading to a maximal  $x_+ = c_{\leftarrow} \in (X, \preceq)$  that is a fixed point of the generating function  $g$  whenever the supremum  $t_{\leftarrow}$  of the chain  $\rightarrow T$  is in  $X$ . ■

*Remark.* The proof of this theorem, despite its apparent length and technically involved character,

carries the highly significant underlying message that

*Any inductive sequential  $g$ -construction of an infinite chained tower  $C_T$  starting with a smallest element  $x_0 \in (X, \preceq)$  such that a supremum  $c_{\leftarrow}$  of the  $g$ -generated sequential chain  $C_T$  in its own tower is contained in itself, must necessarily terminate with a fixed point relation of the type (51) with respect to the supremum. Note from Eqs. (50) and (51) that the role of (ST2) applied to a fully ordered tower is the identification of the maximal of the tower — which depends only on the tower and has nothing to do with anything outside it — with its supremum that depends both on the tower and its complement.*

Thus although purely set-theoretic in nature, the filter-base associated with a sequentially totally ordered set may be interpreted to lead to the usual notions of adherence and convergence of filters and thereby of a generated topology for  $(X, \preceq)$ , see Appendix A.1 and Example A.1.3. This very significant apparent inter-relation between topologies, filters and orderings will form the basis of our approach to the condition of maximal ill-posedness for chaos.

In the second stage of the three-stage programme leading to Zorn’s lemma, the tower Theorem 4.1 and the comments of the preceding paragraph are applied at a higher level to a very special class of the power set of a set, the class of all the chains of a partially ordered set, to directly lead to the physically significant

**Theorem 4.2** (Hausdorff Maximal Principle). Every partially ordered set  $(X, \preceq)$  has a maximal totally ordered subset.<sup>23</sup>

*Proof.* Here the base level is

$$\mathcal{X} = \{C \in \mathcal{P}(X) : C \text{ is a chain in } (X, \preceq)\} \subseteq \mathcal{P}(X) \tag{52}$$

be the set of all the totally ordered subsets of  $(X, \preceq)$ . Since  $\mathcal{X}$  is a collection of (sub)sets of  $X$ , we order it by the inclusion relation on  $\mathcal{X}$  and use the tower Theorem to demonstrate that  $(\mathcal{X}, \subseteq)$  has a maximal element  $C_{\leftarrow}$ , which by the definition of  $\mathcal{X}$ , is the required maximal chain in  $(X, \preceq)$ .

Let  $\mathcal{C}$  be a chain in  $\mathcal{X}$  of the chains in  $(X, \preceq)$ . In order to apply the tower Theorem to  $(\mathcal{X}, \subseteq)$  we need to verify hypothesis (ST2) that the smallest

$$C_* = \sup_{\mathcal{X}} \mathcal{C} = \bigcup_{C \in \mathcal{C}} C \tag{53}$$

of the possible upper bounds of  $\mathcal{C}$  [see Eq. (45)] is a chain of  $(X, \preceq)$ . Indeed, if  $x_1, x_2 \in X$  are two points of  $C_{\text{sup}}$  with  $x_1 \in C_1$  and  $x_2 \in C_2$ , then from the  $\subseteq$ -comparability of  $C_1$  and  $C_2$  we may choose  $x_1, x_2 \in C_1 \supseteq C_2$ , say. Thus  $x_1$  and  $x_2$  are  $\preceq$ -comparable as  $C_1$  is a chain in  $(X, \preceq)$ ;  $C_* \in \mathcal{X}$  is therefore a chain in  $(X, \preceq)$  which establishes that the supremum of a chain of  $(\mathcal{X}, \subseteq)$  is a chain in  $(X, \preceq)$ .

The tower Theorem 4.1 can now be applied to  $(\mathcal{X}, \subseteq)$  with  $C_0$  as its smallest element to construct a  $g$ -sequentially towered fully ordered subset of  $\mathcal{X}$  consisting of chains in  $X$

$$\begin{aligned} \mathcal{C}_T &= \{C_i \in \mathcal{P}(X) : C_i \subseteq C_j \text{ for } i \leq j \in \mathbb{N}\} \\ &= \rightarrow \mathcal{T} \subseteq \mathcal{P}(X) \end{aligned}$$

of  $(\mathcal{X}, \subseteq)$  — consisting of the common elements of all  $g$ -sequential towers  $\mathcal{T} \in \mathfrak{T}$  of  $(\mathcal{X}, \subseteq)$  — that in fact is a principal filter base of chained subsets of  $(X, \preceq)$  at  $C_0$ . The supremum (chain in  $X$ )  $C_{\leftarrow}$  of  $\mathcal{C}_T$  in  $\mathcal{C}_T$  must now satisfy, by Theorem 4.1, the fixed point  $g$ -chain of  $X$

$$\sup_{\mathcal{C}_T} (\mathcal{C}_T) = C_{\leftarrow} = g(C_{\leftarrow}) \in \mathcal{C}_T \subseteq \mathcal{P}(X),$$

where the chain  $g(C) = C \cup f_C(\mathcal{G}(C) - C)$  with  $\mathcal{G}(C) = \{x \in X - C : C \cup \{x\} \in \mathcal{X}\}$ , is an immediate successor of  $C$  obtained by choosing one point  $x = f_C(\mathcal{G}(C) - C)$  from the many possible in  $\mathcal{G}(C) - C$  such that the resulting  $g(C) = C \cup \{x\}$  is a strict successor of the chain  $C$  with no others lying between it and  $C$ . Note that  $C_{\leftarrow} \in (\mathcal{X}, \subseteq)$  is

<sup>23</sup>Recall that this means that if there is a totally ordered chain  $C$  in  $(X, \preceq)$  that succeeds  $C_+$ , then  $C$  must be  $C_+$  so that no chain in  $X$  can be strictly larger than  $C_+$ . The notation adopted here and below is the following: If  $X = \{x, y\}$  is a non-empty set, then  $\mathcal{X} := \mathcal{P}(X) = \{A : A \subseteq X\} = \{\emptyset, \{x\}, \{y\}, \{x, y\}\}$  is the set of subsets of  $X$ , and  $\mathfrak{X} := \mathcal{P}^2(X) = \{\mathcal{A} : \mathcal{A} \subseteq \mathcal{X}\}$ , the set of all subsets of  $\mathcal{X}$ , consists of the 16 elements  $\emptyset, \{\emptyset\}, \{\{x\}\}, \{\{y\}\}, \{\{x, y\}\}, \{\{\emptyset\}, \{x\}\}, \{\{\emptyset\}, \{y\}\}, \{\{\emptyset\}, \{x, y\}\}, \{\{x\}, \{y\}\}, \{\{x\}, \{x, y\}\}, \{\{y\}, \{x, y\}\}, \{\{\emptyset\}, \{x\}, \{y\}\}, \{\{\emptyset\}, \{x\}, \{x, y\}\}, \{\{\emptyset\}, \{y\}, \{x, y\}\}, \{\{x\}, \{y\}, \{x, y\}\}$ , and  $\mathcal{X}$ : an element of  $\mathcal{P}^2(X)$  is a subset of  $\mathcal{P}(X)$ , any element of which is a subset of  $X$ . Thus if  $C = \{0, 1, 2\}$  is a chain in  $(X = \{0, 1, 2\}, \leq)$ , then  $\mathcal{C} = \{\{0\}, \{0, 1\}, \{0, 1, 2\}\} \subseteq \mathcal{P}(X)$  and  $\mathfrak{C} = \{\{\{0\}\}, \{\{0\}, \{0, 1\}\}, \{\{0\}, \{0, 1\}, \{0, 1, 2\}\}\} \subseteq \mathcal{P}^2(X)$  represent chains in  $(\mathcal{P}(X), \subseteq)$  and  $(\mathcal{P}^2(X), \subseteq)$ , respectively.

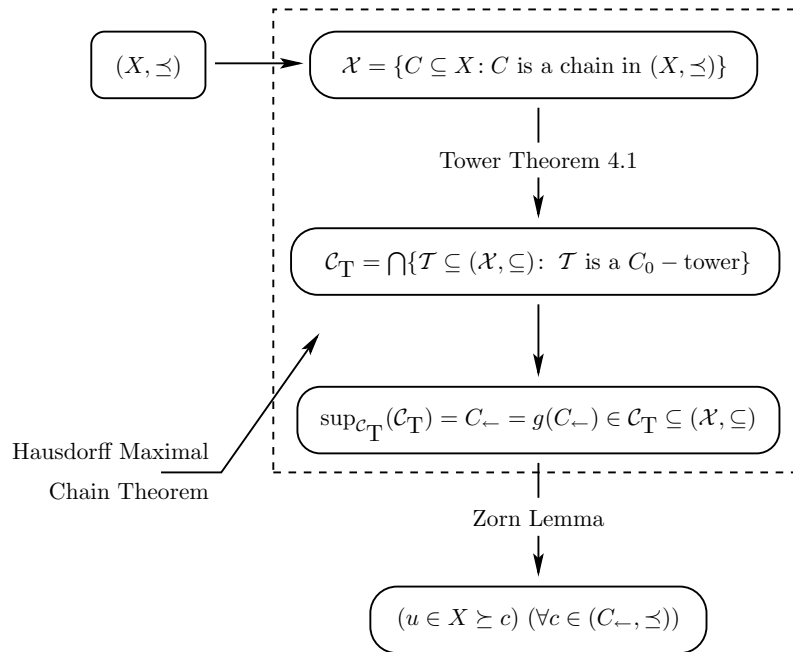


Fig. 10. Application of Zorn’s Lemma to  $(X, \preceq)$ . Starting with a partially ordered set  $(X, \preceq)$ , construct:

- (a) The one-level higher subset  $\mathcal{X} = \{C \in \mathcal{P}(X) : C \text{ is a chain in } (X, \preceq)\}$  of  $\mathcal{P}(X)$  consisting of all the totally ordered subsets of  $(X, \preceq)$ ,
- (b) The smallest common  $g$ -sequential totally ordered towered chain  $\mathcal{C}_T = \{C_i \in \mathcal{P}(X) : C_i \subseteq C_j \text{ for } i \leq j\} \subseteq \mathcal{P}(X)$  of all sequential  $g$ -towers of  $\mathcal{X}$  by Theorem 4.1, which in fact is a principal filter base of totally ordered subsets of  $(X, \preceq)$  at the smallest element  $C_0$ .
- (c) Apply Hausdorff Maximal Principle to  $(\mathcal{X}, \subseteq)$  to get the subset  $\sup_{\mathcal{C}_T}(\mathcal{C}_T) = C_{\leftarrow} = g(C_{\leftarrow}) \in \mathcal{C}_T \subseteq \mathcal{P}(X)$  of  $(X, \preceq)$  as the supremum of  $(\mathcal{X}, \subseteq)$  in  $\mathcal{C}_T$ . The identification of this supremum as a maximal element of  $(\mathcal{X}, \subseteq)$  is a consequence of (ST2) and Eqs. (50), (51) that actually puts the supremum into  $\mathcal{X}$  itself.

By returning to the original level  $(X, \preceq)$

- (d) Zorn’s Lemma finally yields the required maximal element  $u \in X$  as an upper bound of the maximal totally ordered subset  $(C_{\leftarrow}, \preceq)$  of  $(X, \preceq)$ .

The dashed segment denotes the higher Hausdorff  $(\mathcal{X}, \subseteq)$  level leading to the base  $(X, \preceq)$  Zorn level.

only one of the many maximal fully ordered subsets possible in  $(X, \preceq)$ . ■

With the assurance of the existence of a maximal chain  $C_{\leftarrow}$  among all fully ordered subsets of a partially ordered set  $(X, \preceq)$ , the arguments are completed by returning to the basic level of  $X$ .

**Theorem 4.3** (Zorn’s Lemma). *Let  $(X, \preceq)$  be a partially ordered set such that every totally ordered subset of  $X$  has an upper bound in  $X$ . Then  $X$  has at least one maximal element with respect to its order.*

*Proof.* The proof of this final part is a mere application of the Hausdorff Maximal Principle on the existence of a maximal chain  $C_{\leftarrow}$  in  $X$  to the hypothesis of this theorem that  $C_{\leftarrow}$  has an upper bound  $u$  in  $X$  that quickly leads to the identification of this bound as a maximal element  $x_+$  of  $X$ .

Indeed, if there is an element  $v \in X$  that is comparable to  $u$  and  $v \not\preceq u$ , then  $v$  cannot be in  $C_{\leftarrow}$  as it is necessary for every  $x \in C_{\leftarrow}$  to satisfy  $x \preceq u$ . Clearly then  $C_{\leftarrow} \cup \{v\}$  is a chain in  $(X, \preceq)$  bigger than  $C_{\leftarrow}$  which contradicts the assumed maximality of  $C_{\leftarrow}$  among the chains of  $X$ . ■

The sequence of steps leading to Zorn’s Lemma, and thence to the maximal of a partially ordered set, is summarized in Fig. 10.

The three examples below of the application of Zorn’s Lemma clearly reflect the increasing complexity of the problem considered, with the maximal a point, a subset, and a set of subsets of  $X$ , so that these are elements of  $X$ ,  $\mathcal{P}(X)$  and  $\mathcal{P}^2(X)$ , respectively.

**Example 4.2**

- (1) Let  $X = (\{a, b, c\}, \preceq)$  be a three-point base-

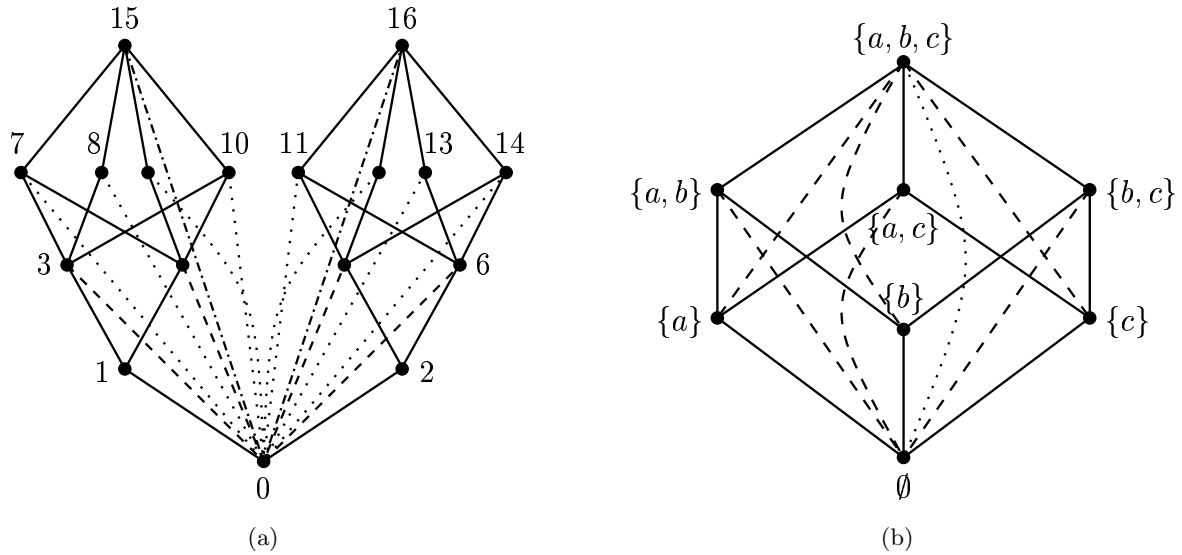


Fig. 11. Tree diagrams of two partially ordered sets where two points are connected by a line iff they are comparable to each other, with the solid lines linking immediate neighbors and the dashed, dotted and dashed-dotted lines denoting second, third and fourth generation orderings according to the principle of transitivity of the order relation. There are  $8 \times 2$  chains of (a) and 7 chains of (b) starting from respective smallest elements with the immediate successor chains shown in solid lines. The 17 point set  $X = \{0, 1, 2, \dots, 15, 16\}$  in (a) has two maximals but no maximum, while in (b) there is a single maximum of  $\mathcal{P}(\{a, b, c\})$ , and three maximals without any maximum for  $\mathcal{P}(\{a, b, c\}) - \{a, b, c\}$ . In (a), let  $A = \{1, 3, 4, 7, 9, 10, 15\}$ ,  $B = \{1, 3, 4, 6, 7, 13, 15\}$ ,  $C = \{1, 3, 4, 10, 11, 16\}$  and  $D = \{1, 3, 4\}$ . The upper bounds of  $D$  in  $A$  are 7, 10 and 15 without any supremum (as there is no smallest element of  $\{7, 10, 15\}$ ), and the upper bounds of  $D$  in  $B$  are 7 and 15 with  $\sup_B(D) = 7$ , while  $\sup_C(D) = 10$ . Finally the maximal, maximum and the supremum in  $A$  of  $\{1, 3, 4, 7\}$  are all the same illustrating how the supremum of a set can belong to itself. Observe how the supremum and upper bound of a set are with reference to its complement in contrast with the maximum and maximal that have nothing to do with anything outside the set.

level ground set ordered lexicographically, that is  $a \prec b \prec c$ . A chain  $\mathcal{C}$  of the partially ordered Hausdorff-level set  $\mathcal{X}$  consisting of subsets of  $X$  given by Eq. (52) is, for example,  $\{\{a\}, \{a, b\}\}$  and the six  $g$ -sequential chained towers

$$\begin{aligned} \mathcal{C}_1 &= \{\emptyset, \{a\}, \{a, b\}, \{a, b, c\}\}, \\ \mathcal{C}_2 &= \{\emptyset, \{a\}, \{a, c\}, \{a, b, c\}\} \\ \mathcal{C}_3 &= \{\emptyset, \{b\}, \{a, b\}, \{a, b, c\}\}, \\ \mathcal{C}_4 &= \{\emptyset, \{b\}, \{b, c\}, \{a, b, c\}\} \\ \mathcal{C}_5 &= \{\emptyset, \{c\}, \{a, c\}, \{a, b, c\}\}, \\ \mathcal{C}_6 &= \{\emptyset, \{c\}, \{b, c\}, \{a, b, c\}\} \end{aligned}$$

built from the smallest element  $\emptyset$  corresponding to the six distinct ways of reaching  $\{a, b, c\}$  from  $\emptyset$  along the sides of the cube marked on the figure with solid lines, all belong to  $\mathcal{X}$ ; see Fig. 11(b). An example of a tower in  $(\mathcal{X}, \subseteq)$  which is not a chain is  $\mathcal{T} = \{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, \{a, b, c\}\}$ . Hence the common infimum towered chained subset

$$\mathcal{C}_T = \{\emptyset, \{a, b, c\}\} \Rightarrow \mathcal{T} \subseteq \mathcal{P}(X)$$

of  $\mathcal{X}$ , with

$$\sup_{\mathcal{C}_T}(\mathcal{C}_T) = C_{\leftarrow} = \{a, b, c\} = g(C_{\leftarrow}) \in \mathcal{C}_T \subseteq \mathcal{P}(X)$$

the only maximal element of  $\mathcal{P}(X)$ . Zorn's Lemma now assures the existence of a maximal element of  $c \in X$ . Observe how the maximal element of  $(X, \preceq)$  is obtained by going one level higher to  $\mathcal{X}$  at the Hausdorff stage and returning to the base level  $X$  at Zorn, see Fig. 10 for a schematic summary of this sequence of steps.

(2) *Basis of a vector space.* A linearly independent set of vectors in a vector space  $X$  that spans the space is known as the Hamel basis of  $X$ . To prove the existence of a Hamel basis in a vector space, Zorn's lemma is invoked as follows.

The ground base level of the linearly independent subsets of  $X$

$$\begin{aligned} \mathcal{X} &= \{\{x_{i_j}\}_{j=1}^J \in \mathcal{P}(X) : \text{Span}(\{x_{i_j}\}_{j=1}^J) \\ &= 0 \Rightarrow (\alpha_j)_{j=1}^J = 0 \forall J \geq 1\} \subseteq \mathcal{P}(X), \end{aligned}$$

with  $\text{Span}(\{x_{i_j}\}_{j=1}^J) := \sum_{j=1}^J \alpha_j x_{i_j}$ , is such that no  $x \in \mathcal{X}$  can be expressed as a linear combination of

the elements of  $\mathcal{X} - \{x\}$ .  $\mathcal{X}$  clearly has a smallest element, say  $\{x_{i_1}\}$ , for some non-zero  $x_{i_1} \in X$ . Let the higher Hausdorff level

$\mathfrak{X} = \{\mathcal{C} \in \mathcal{P}^2(X) : \mathcal{C} \text{ is a chain in } (\mathcal{X}, \subseteq)\} \subseteq \mathcal{P}^2(X)$   
and collection of the chains

$$\mathcal{C}_{i_K} = \{\{x_{i_1}\}, \{x_{i_1}, x_{i_2}\}, \dots, \{x_{i_1}, x_{i_2}, \dots, x_{i_K}\}\} \in \mathcal{P}^2(X)$$

of  $\mathcal{X}$  comprising linearly independent subsets of  $X$  be  $g$ -built from the smallest  $\{x_{i_1}\}$ . Any chain  $\mathcal{C}$  of  $\mathfrak{X}$  is bounded above by the union  $\mathcal{C}_* = \sup_{\mathfrak{X}} \mathcal{C} = \bigcup_{\mathcal{C} \in \mathfrak{X}} \mathcal{C}$  which is a chain in  $\mathcal{X}$  containing  $\{x_{i_1}\}$ , thereby verifying (ST2) for  $\mathfrak{X}$ . Application of the tower theorem to  $\mathfrak{X}$  implies that the element

$$\mathfrak{C}_T = \{\mathcal{C}_{i_1}, \mathcal{C}_{i_2}, \dots, \mathcal{C}_{i_n}, \dots\} \Rightarrow \mathfrak{X} \subseteq \mathcal{P}^2(X)$$

in  $\mathfrak{X}$  of chains of  $\mathcal{X}$  is a  $g$ -sequential fully ordered towered subset of  $(\mathfrak{X}, \subseteq)$  consisting of the common elements of all  $g$ -sequential towers of  $(\mathfrak{X}, \subseteq)$ , that in fact is a *chained principal ultrafilter on  $(\mathcal{P}(X), \subseteq)$  generated by the filter-base  $\{\{x_{i_1}\}\}$  at  $\{x_{i_1}\}$* , where

$$\mathfrak{X} = \{\mathcal{C}_{i_1}, \mathcal{C}_{i_2}, \dots, \mathcal{C}_{j_n}, \mathcal{C}_{j_{n+1}}, \dots\}$$

for some  $n \in \mathbb{N}$  is an example of non-chained  $g$ -tower whenever  $(\mathcal{C}_{j_k})_{k=n}^\infty$  is neither contained in nor contains any member of the  $(\mathcal{C}_{i_k})_{k=1}^\infty$  chain. Hausdorff's chain theorem now yields the fixed-point  $g$ -chain  $\mathcal{C}_\leftarrow \in \mathfrak{X}$  of  $\mathcal{X}$

$$\begin{aligned} \sup_{\mathfrak{C}_T}(\mathfrak{C}_T) &= \mathcal{C}_\leftarrow = \{\{x_{i_1}\}, \{x_{i_1}, x_{i_2}\}, \{x_{i_1}, x_{i_2}, x_{i_3}\}, \dots\} \\ &= g(\mathcal{C}_\leftarrow) \in \mathfrak{C}_T \subseteq \mathcal{P}^2(X) \end{aligned}$$

as a maximal *totally ordered principal filter on  $X$  that is generated by the filter-base  $\{\{x_{i_1}\}\}$  at  $x_{i_1}$* , whose supremum  $B = \{x_{i_1}, x_{i_2}, \dots\} \in \mathcal{P}(X)$  is, by Zorn's lemma, a maximal element of the base level  $\mathcal{X}$ . This maximal linearly independent subset of  $X$  is the required Hamel basis for  $X$ : Indeed, if the span of  $B$  is not the whole of  $X$ , then  $\text{Span}(B) \cup x$ , with  $x \notin \text{Span}(B)$  would, by definition, be a linearly independent set of  $X$  strictly larger than  $B$ , contradicting the assumed maximality of the later. It needs to be understood that since the infinite basis cannot be classified as being linearly independent, we have here an important example of the supremum of the maximal chained set not belonging to the set even though this criterion was explicitly used in the construction process according to (ST2) and (ST3).

Compared to this purely algebraic concept of basis in a vector space, is the Schauder basis in a normed space which combines topological structure with the linear in the form of convergence: If a normed vector space contains a sequence  $(e_i)_{i \in \mathbb{Z}_+}$  with the property that for every  $x \in X$  there is a unique sequence of scalars  $(\alpha_i)_{i \in \mathbb{Z}_+}$  such that the remainder  $\|x - (\alpha_1 e_1 + \alpha_2 e_2 + \dots + \alpha_I e_I)\|$  approaches 0 as  $I \rightarrow \infty$ , then the collection  $(e_i)$  is known as a Schauder basis for  $X$ .

(3) *Ultrafilter.* Let  $X$  be a set. The set  ${}_F\mathcal{S} = \{S_\alpha \in \mathcal{P}(X) : S_\alpha \cap S_\beta \neq \emptyset, \forall \alpha \neq \beta\} \subseteq \mathcal{P}(X)$  of all nonempty subsets of  $X$  with finite intersection property is known as a *filter subbase on  $X$*  and  ${}_F\mathcal{B} = \{B \subseteq X : B = \bigcap_{i \in I} S_i\}$ , for  $I \subset \mathbb{D}$  a finite subset of a directed set  $\mathbb{D}$ , is a *filter-base on  $X$  associated with the subbase  ${}_F\mathcal{S}$* ; cf. Appendix A.1. Then the *filter generated by  ${}_F\mathcal{S}$*  consisting of every superset of the finite intersections  $B \in {}_F\mathcal{B}$  of sets of  ${}_F\mathcal{S}$  is the smallest filter that contains the subbase  ${}_F\mathcal{S}$  and base  ${}_F\mathcal{B}$ . For notational simplicity, we will denote the subbase  ${}_F\mathcal{S}$  in the rest of this example simply by  $\mathcal{S}$ .

Consider the base-level ground set of all filter subbases on  $X$

$$\begin{aligned} \mathfrak{S} &= \left\{ \mathcal{S} \in \mathcal{P}^2(X) : \bigcap_{\emptyset \neq \mathcal{R} \subseteq \mathcal{S}} \mathcal{R} \neq \emptyset \text{ for every finite subset of } \mathcal{S} \right\} \\ &\subseteq \mathcal{P}^2(X), \end{aligned}$$

ordered by inclusion in the sense that  $\mathcal{S}_\alpha \subseteq \mathcal{S}_\beta$  for all  $\alpha \preceq \beta \in \mathbb{D}$ , and let the higher Hausdorff-level

$$\tilde{\mathfrak{X}} = \{\mathfrak{C} \in \mathcal{P}^3(X) : \mathfrak{C} \text{ is a chain in } (\mathfrak{S}, \subseteq)\} \subseteq \mathcal{P}^3(X)$$

comprising the collection of the totally ordered chains

$$\mathfrak{C}_\kappa = \{\{S_\alpha\}, \{S_\alpha, S_\beta\}, \dots, \{S_\alpha, S_\beta, \dots, S_\kappa\}\} \in \mathcal{P}^3(X)$$

of  $\mathfrak{S}$  be  $g$ -built from the smallest  $\{S_\alpha\}$  then an *ultrafilter on  $X$*  is a maximal member  $\mathcal{S}_+$  of  $(\mathfrak{S}, \subseteq)$  in the usual sense that any subbase  $\mathcal{S}$  on  $X$  must necessarily be contained in  $\mathcal{S}_+$  so that  $\mathcal{S}_+ \subseteq \mathcal{S} \Rightarrow \mathcal{S} = \mathcal{S}_+$  for any  $\mathcal{S} \subseteq \mathcal{P}(X)$  with FIP. The tower theorem now implies that the element

$$\tilde{\mathfrak{C}}_T = \{\mathfrak{C}_\alpha, \mathfrak{C}_\beta, \dots, \mathfrak{C}_\nu, \dots\} = \tilde{\mathfrak{X}} \subseteq \mathcal{P}^3(X)$$

of  $\mathcal{P}^4(X)$ , which is a chain in  $\tilde{\mathfrak{X}}$  of the chains of  $\mathfrak{S}$ , is a  $g$ -sequential fully ordered towered subset of the common elements of all sequential towers of  $(\tilde{\mathfrak{X}}, \subseteq)$

that is a *chained principal ultrafilter* on  $(\mathcal{P}^2(X), \subseteq)$  generated by the filter-base  $\{\{S_\alpha\}\}$  at  $\{S_\alpha\}$ , where

$$\tilde{\mathfrak{F}} = \{\mathfrak{C}_\alpha, \mathfrak{C}_\beta, \dots, \mathfrak{C}_\sigma, \mathfrak{C}_\tau, \dots\},$$

is an obvious example of non-chained  $g$ -tower whenever  $(\mathfrak{C}_\sigma)$  is neither contained in, nor contains, any member of the  $\mathfrak{C}_\alpha$ -chain. Hausdorff's chain theorem now yields the fixed-point  $\mathfrak{C}_{\leftarrow} \in \tilde{\mathfrak{F}}$

$$\begin{aligned} \sup_{\mathfrak{C}_T}(\tilde{\mathfrak{C}}_T) &= \mathfrak{C}_{\leftarrow} = \{\{S_\alpha\}, \{S_\alpha, S_\beta\}, \{S_\alpha, S_\beta, S_\gamma\}, \dots\} \\ &= g(\mathfrak{C}_{\leftarrow}) \in \tilde{\mathfrak{C}}_T \subseteq \mathcal{P}^3(X) \end{aligned}$$

as a maximal *totally ordered*  $g$ -chained towered subset of  $X$  that is, by Zorn's lemma, a maximal element of the base level subset  $\mathfrak{S}$  of  $\mathcal{P}^2(X)$ .  $\mathfrak{C}_{\leftarrow}$  is a *chained principal ultrafilter* on  $(\mathcal{P}(X), \subseteq)$  generated by the filter-base  $\{\{S_\alpha\}\}$  at  $S_\alpha$ , while  $\mathcal{S}_+ = \{S_\alpha, S_\beta, S_\gamma, \dots\} \in \mathcal{P}^2(X)$  is an (non-principal) *ultrafilter* on  $X$  — characterized by the property that any collection of subsets on  $X$  with FIP (that is any filter subbase on  $X$ ) must be contained in the maximal set  $\mathcal{S}_+$  having FIP — that is not a principal filter unless  $\mathcal{S}_\alpha$  is a singleton set  $\{x_\alpha\}$ .

What emerges from these applications of Zorn's Lemma is the remarkable fact that *infinities (the dot-dot-dots) can be formally introduced as "limiting cases" of finite systems in a purely set-theoretic context without the need for topologies, metrics or convergences*. The significance of this observation will become clear from our discussions on filters and topology leading to Sec. 4.2 below. Also, the observation on the successive iterates of the power sets  $\mathcal{P}(X)$  in the examples above was to suggest their anticipated role in the complex evolution of a dynamical system that is expected to play a significant part in our future interpretation and understanding of this adaptive and self-organizing phenomenon of nature.

**End Tutorial 5**

---

From the examples in Tutorial 5, it should be clear that the sequential steps summarized in Fig. 10 are involved in an application of Zorn's lemma to show that a partially ordered set has a maximal element with respect to its order. Thus for a partially ordered set  $(X, \preceq)$ , form the set  $\mathcal{X}$  of all chains  $C$  in  $X$ . If  $C_+$  is a maximal chain of  $X$  obtained by the Hausdorff Maximal Principle from the chain  $\mathcal{C}$  of all chains of  $X$ , then its supremum  $u$  is a maximal

element of  $(X, \preceq)$ . This sequence is now applied, as in Example 4.2(1), to the set of arbitrary relations  $\text{Multi}(X)$  on an infinite set  $X$  in order to formulate our definition of chaos that follows.

Let  $f$  be a *noninjective map* in  $\text{Multi}(X)$  and  $P(f)$  the number of injective branches of  $f$ . Denote by

$$F = \{f \in \text{Multi}(X) : f \text{ is a noninjective function on } X\} \subseteq \text{Multi}(X)$$

the resulting basic collection of noninjective functions in  $\text{Multi}(X)$ .

- (i) For every  $\alpha$  in some directed set  $\mathbb{D}$ , let  $F$  have the extension property

$$(\forall f_\alpha \in F)(\exists f_\beta \in F) : P(f_\alpha) \leq P(f_\beta)$$

- (ii) Let a partial order  $\preceq$  on  $\text{Multi}(X)$  be defined, for  $f_\alpha, f_\beta \in \text{Map}(X) \subseteq \text{Multi}(X)$  by

$$P(f_\alpha) \leq P(f_\beta) \Leftrightarrow f_\alpha \preceq f_\beta, \tag{54}$$

with  $P(f) := 1$  for the smallest  $f$ , define a partially ordered subset  $(F, \preceq)$  of  $\text{Multi}(X)$ . This is actually a preorder on  $\text{Multi}(X)$  in which functions with the same number of injective branches are equivalent to each other.

- (iii) Let

$$\begin{aligned} C_\nu &= \{f_\alpha \in \text{Multi}(X) : f_\alpha \preceq f_\nu\} \in \mathcal{P}(F), \\ &\nu \in \mathbb{D}, \end{aligned}$$

be  $g$ -chains of non-injective functions of  $\text{Multi}(X)$  and

$$\mathcal{X} = \{C \in \mathcal{P}(F) : C \text{ is a chain in } (F, \preceq)\} \subseteq \mathcal{P}(F)$$

denote the corresponding Hausdorff level of all chains of  $F$ , with

$$\mathcal{C}_T = \{C_\alpha, C_\beta, \dots, C_\nu, \dots\} \Rightarrow T \subseteq \mathcal{P}(F)$$

being a  $g$ -sequential in  $\mathcal{X}$ . By Hausdorff Maximal Principle, there is a maximal fixed-point  $g$ -towered chain  $C_{\leftarrow} \in \mathcal{X}$  of  $F$

$$\begin{aligned} \sup_{\mathcal{C}_T}(C_T) &= C_{\leftarrow} = \{f_\alpha, f_\beta, \dots\} \\ &= g(C_{\leftarrow}) \in \mathcal{C}_T \subseteq \mathcal{P}(F). \end{aligned}$$

Zorn's Lemma applied to this maximal chain yields its supremum as the maximal element of  $C_{\leftarrow}$ , and thereby of  $F$ . It needs to be appreciated, as in the case of the algebraic Hamel basis, that the existence of this maximal non-functional element was

obtained purely set theoretically as the “limit” of a net of functions with increasing nonlinearity, without resorting to any topological arguments. Because it is not a function, this supremum does not belong to the functional  $g$ -towered chain having it as a fixed point, and this maximal chain does not possess a largest, or even a maximal, element, although it does have a supremum.<sup>24</sup> The supremum is a contribution of the inverse functional relations ( $f_\alpha^-$ ) in the following sense. From Eq. (2), the net of increasingly non-injective functions of Eq. (54) implies a corresponding net of increasingly multi-valued functions ordered inversely by the inverse relation  $f_\alpha \preceq f_\beta \Leftrightarrow f_\beta^- \preceq f_\alpha^-$ . Thus the inverse relations which are as much an integral part of graphical convergence as are the direct relations, have a smallest element belonging to the multifunctional class. Clearly, this smallest element as the required supremum of the increasingly non-injective tower of functions defined by Eq. (54), serves to complete the significance of the tower by capping it with a “boundary” element that can be taken to bridge the classes of functional and non-functional relations on  $X$ .

We are now ready to define a *maximally ill-posed problem*  $f(x) = y$  for  $x, y \in X$  in terms of a *maximally non-injective map*  $f$  as follows.

**Definition 4.1** (Chaotic map). Let  $A$  be a non-empty closed set of a compact Hausdorff space  $X$ . A function  $f \in \text{Multi}(X)$  equivalently the sequence of functions ( $f_i$ ) is maximally non-injective or chaotic on  $A$  with respect to the order relation (54) if

- (a) for any  $f_i$  on  $A$  there exists an  $f_j$  on  $A$  satisfying  $f_i \preceq f_j$  for every  $j > i \in \mathbb{N}$ .
- (b) the set  $\mathcal{D}_+$  consists of a countable collection of isolated singletons.

**Definition 4.2** (Maximally ill-posed problem). Let  $A$  be a non-empty closed set of a compact Hausdorff space  $X$  and let  $f$  be a functional relation in  $\text{Multi}(X)$ . The problem  $f(x) = y$  is maximally ill-posed at  $y$  if  $f$  is chaotic on  $A$ .

As an example of the application of these definitions, on the dense set  $\mathcal{D}_+$ , the tent map satisfies both the conditions of sensitive dependence on initial conditions and topological transitivity

[Devaney, 1989] and is also maximally non-injective; the tent map is therefore chaotic on  $\mathcal{D}_+$ . In contrast, the examples of Secs. 1 and 2 are not chaotic as the maps are not topologically transitive, although the Liapunov exponents, as in the case of the tent map, are positive. Here the ( $f_n$ ) are identified with the iterates of  $f$ , and the “fixed point” as one through which graphs of all the functions on residual index subsets pass. When the set of points  $\mathcal{D}_+$  is dense in  $[0, 1]$  and both  $\mathcal{D}_+$  and  $[0, 1] - \mathcal{D}_+ = [0, 1] - \bigcup_{i=0}^{\infty} f^{-i}(\text{Per}(f))$  (where  $\text{Per}(f)$  denotes the set of periodic points of  $f$ ) are totally disconnected, it is expected that at any point on this complement the behavior of the limit will be similar to that on  $\mathcal{D}_+$ : these points are special as they tie up the iterates on  $\text{Per}(f)$  to yield the multifunctions. Therefore in any neighborhood  $U$  of a  $\mathcal{D}_+$ -point, there is an  $x_0$  at which the *forward orbit*  $\{f^i(x_0)\}_{i \geq 0}$  is *chaotic* in the sense that

- (a) the sequence neither diverges nor does it converge in the image space of  $f$  to a periodic orbit of any period, and
- (b) the Liapunov exponent is given by

$$\begin{aligned} \lambda(x_0) &\stackrel{\text{def}}{=} \lim_{n \rightarrow \infty} \ln \left| \frac{df^n(x_0)}{dx} \right|^{1/n} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} \ln \left| \frac{df(x_i)}{dx} \right|, \quad x_i = f^i(x_0), \end{aligned}$$

which is a measure of the average slope of an orbit at  $x_0$  or equivalently of the average loss of information of the position of a point after one iteration, is positive. Thus *an orbit with positive Liapunov exponent is chaotic if it is not asymptotic* (that is neither convergent nor adherent, having no convergent sub-orbit in the sense of Appendix A.1) *to an unstable periodic orbit or to any other limit set on which the dynamics is simple*. A basic example of a chaotic orbit is that of an irrational in  $[0, 1]$  under the shift map and that of the chaotic set its closure, the full unit interval.

Let  $f \in \text{Map}((X, \mathcal{U}))$  and suppose that  $A = \{f^j(x_0)\}_{j \in \mathbb{N}}$  is a sequential set corresponding to the orbit  $\text{Orb}(x_0) = (f^j(x_0))_{j \in \mathbb{N}}$ , and let  $f_{\mathbb{R}_i}(x_0) = \bigcup_{j \geq i} f^j(x_0)$  be the  $i$ -residual of the sequence  $(f^j(x_0))_{j \in \mathbb{N}}$ , with  ${}_{\mathbb{F}}\mathcal{B}_{x_0} = \{f_{\mathbb{R}_i}(x_0) : \text{Res}(\mathbb{N}) \rightarrow X$

<sup>24</sup>A similar situation arises in the following more intuitive example. Although the subset  $A = \{1/n\}_{n \in \mathbb{Z}_+}$  of the interval  $I = [-1, 1]$  has no smallest or minimal elements, it does have the infimum 0. Likewise, although  $A$  is bounded below by any element of  $[-1, 0)$ , it has no greatest lower bound in  $[-1, 0) \cup (0, 1]$ .



for all  $i \in \mathbb{N}$  being the decreasingly nested filter-base associated with  $\text{Orb}(x_0)$ . The so-called  $\omega$ -limit set of  $x_0$  given by

$$\begin{aligned} \omega(x_0) &\stackrel{\text{def}}{=} \{x \in X : (\exists n_k \in \mathbb{N})(n_k \rightarrow \infty)(f^{n_k}(x_0) \rightarrow x)\} \\ &= \{x \in X : (\forall N \in \mathcal{N}_x)(\forall f_{\mathbb{R}_i} \in_F \mathcal{B}_{x_0}) \\ &\quad (f_{\mathbb{R}_i}(x_0) \cap N \neq \emptyset)\} \end{aligned} \tag{55}$$

is simply the adherence set  $\text{adh}(f^j(x_0))$  of the sequence  $(f^j(x_0))_{j \in \mathbb{N}}$ , see Eq. (A.39); hence Definition A.1.11 of the filter-base associated with a sequence and Eqs. (A.16), (A.24), (A.31) and (A.34) allow us to express  $\omega(x_0)$  more meaningfully as

$$\omega(x_0) = \bigcap_{i \in \mathbb{N}} \text{Cl}(f_{\mathbb{R}_i}(x_0)). \tag{56}$$

It is clear from the second of Eqs. 55) that for a continuous  $f$  and any  $x \in X$ ,  $x \in \omega(x_0)$  implies  $f(x) \in \omega(x_0)$  so that the entire orbit of  $x$  lies in  $\omega(x_0)$  whenever  $x$  does imply that the  $\omega$ -limit set is positively invariant; it is also closed because the adherent set is a closed set according to Theorem A.1.3. Hence  $x_0 \in \omega(x_0) \Rightarrow A \subseteq \omega(x_0)$  reduces the  $\omega$ -limit set to the closure of  $A$  without any isolated points,  $A \subseteq \text{Der}(A)$ . In terms of Eq. (A.33) involving principal filters, Eq. (56) in this case may be expressed in the more transparent form  $\omega(x_0) = \bigcap \text{Cl}_F(\mathcal{P}(\{f^j(x_0)\}_{j=0}^\infty))$  where the principal filter  ${}_F\mathcal{P}(\{f^j(x_0)\}_{j=0}^\infty)$  at  $A$  consists of all supersets of  $A = \{f^j(x_0)\}_{j=0}^\infty$ , and  $\omega(x_0)$  represents the adherence set of the principal filter at  $A$ , see the discussion following Theorem A.1.3. If  $A$  represents a chaotic orbit under this condition, then  $\omega(x_0)$  is sometimes known as a *chaotic set* [Alligood *et al.*, 1997]; thus the chaotic orbit infinitely often visits every member of its chaotic set<sup>25</sup> which is simply the  $\omega$ -limit set of a chaotic orbit that is itself contained in its own limit set. Clearly the chaotic set is positive invariant, and from Theorem A.1.3 and its corollary it is also compact. Furthermore, if all (sub)sequences emanating from points  $x_0$  in some neighborhood of the set converge to it, then  $\omega(x_0)$  is called a *chaotic attractor*, see [Alligood *et al.*, 1997]. As common examples of chaotic sets that are not attractors mention may be made of the tent map with a peak value larger than 1 at 0.5, and the logistic map with  $\lambda \geq 4$  again with a peak value at 0.5 exceeding 1.

It is important that the difference in the dynamical behavior of the system on  $\mathcal{D}_+$  and its complement be appreciated. At any fixed point  $x$  of  $f^i$  in  $\mathcal{D}_+$  (or at its equivalent images in  $[x]$ ) the dynamics eventually gets attached to the (equivalent) fixed point, and the sequence of iterates converges graphically in  $\text{Multi}(X)$  to  $x$  (or its equivalent points). When  $x \notin \mathcal{D}_+$ , however, the orbit  $A = \{f^i(x)\}_{i \in \mathbb{N}}$  is chaotic in the sense that  $(f^i(x))$  is not asymptotically periodic and not being attached to any particular point they wander about in the closed chaotic set  $\omega(x) = \text{Der}(A)$  containing  $A$  such that for any given point in the set, some subsequence of the chaotic orbit gets arbitrarily close to it. Such sequences do not converge anywhere but only frequent every point of  $\text{Der}(A)$ . Thus although in the limit of progressively larger iterations there is complete uncertainty of the outcome of an experiment conducted at either of these two categories of initial points, whereas on  $\mathcal{D}_+$  this is due to a random choice from a multifunctional set of equally probable outputs as dictated by the specific conditions under which the experiment was conducted at that instant, on its complement the uncertainty is due to the chaotic behavior of the functional iterates themselves. Nevertheless it must be clearly understood that *this later behavior is entirely due to the multifunctional limits at the  $\mathcal{D}_+$  points which completely determine the behavior of the system on its complement*. As an explicit illustration of this situation, recall that for the shift map  $2x \bmod(1)$  the  $\mathcal{D}_+$  points are the rationals on  $[0, 1]$ , and any irrational is represented by a non-terminating and non-repeating decimal so that almost all decimals in  $[0, 1]$  in any base contain all possible sequences of any number of digits. For the logistic map, the situation is more complex, however. Here the onset of chaos marking the end of the period doubling sequence at  $\lambda_* = 3.5699456$  is signaled by the disappearance of all stable fixed points, Fig. 13(c), with Fig. 13(a) being a demonstration of the stable limits for  $\lambda = 3.569$  that show up as convergence of the iterates to constant valued functions (rather than as constant valued inverse functions) at stable fixed points, shown more emphatically in Fig. 12(a). What actually happens at  $\lambda_*$  is shown in Fig. 16(a) in the next subsection: the almost vertical lines produced at a large, but finite, iterations  $i$

<sup>25</sup>How does this happen for  $A = \{f^i(x_0)\}_{i \in \mathbb{N}}$  that is not the constant sequence  $(x_0)$  at a fixed point? As  $i \in \mathbb{N}$  increases, points are added to  $\{x_0, f(x_0), \dots, f^I(x_0)\}$  not, as would be the case in a normal sequence, as a piled-up Cauchy tail, but as points generally lying between those already present; recall a typical graph as of Fig. 9, for example.

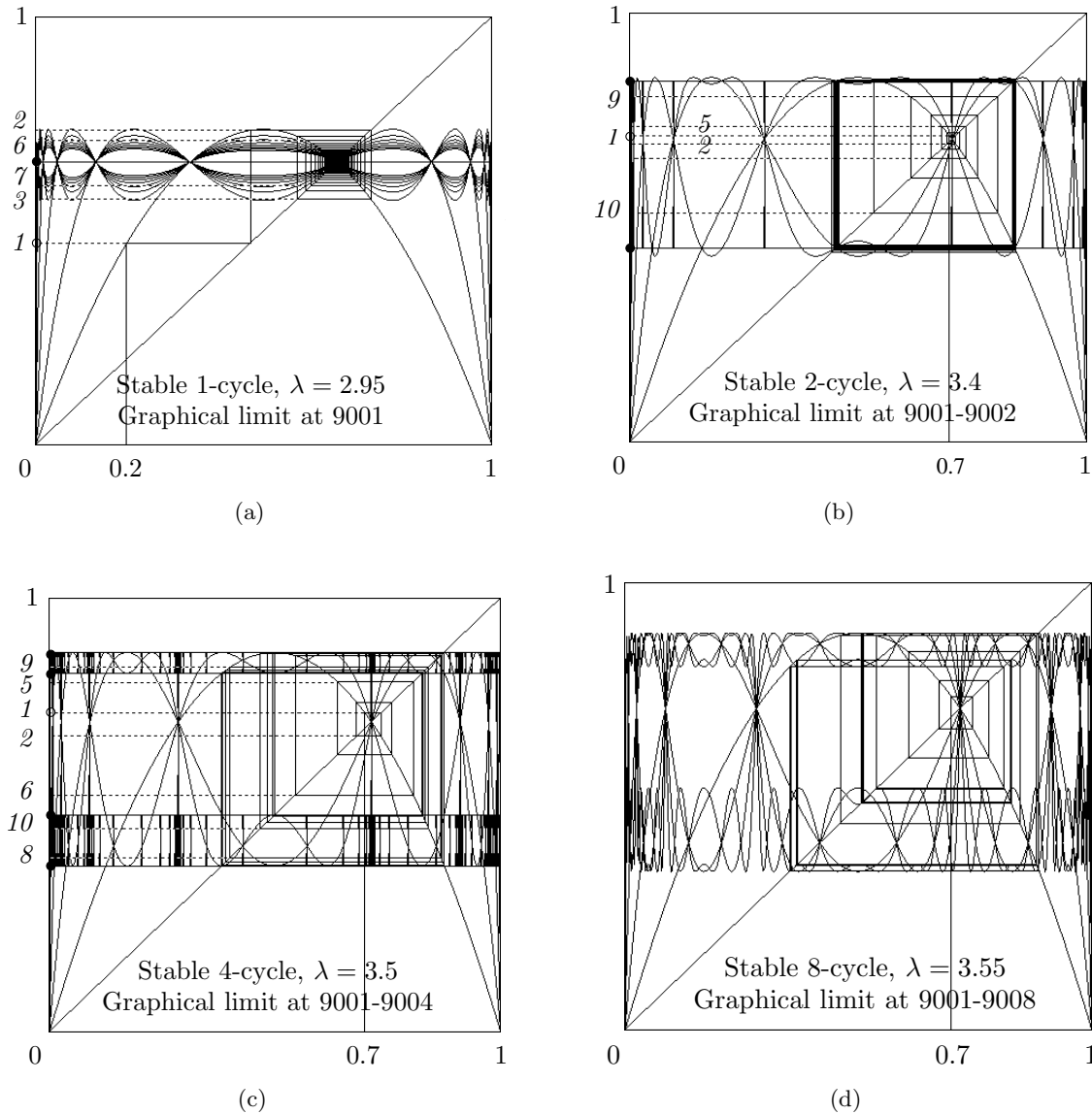
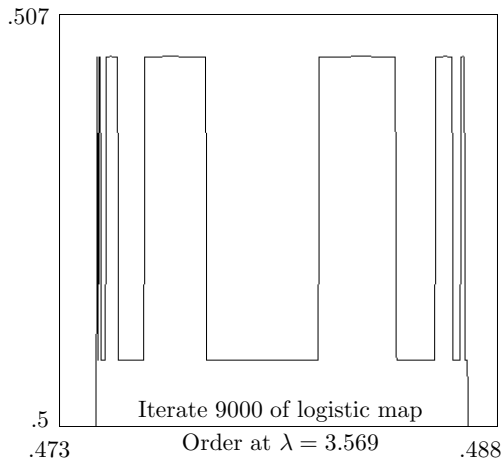


Fig. 12. Fixed points and cycles of logistic map. The isolated fixed point of (b) yields two non-fixed points to which the iterates converge *simultaneously* in the sense that the generated sequence converges to one iff it converges to the other. This suggests that nonlinear dynamics of a system can lead to a situation in which sequences in a Hausdorff space may converge to more than one point. Since convergence depends on the topology (Corollary to Theorem A.1.5), this may be interpreted to mean that nonlinearity tends to modify the basic structure of a space. The sequence of points generated by the iterates of the map are marked on the  $y$ -axis of (a)–(c) in *italics*. The singletons  $\{x\}$  are  $\omega$ -limit sets of the respective fixed point  $x$  and is generated by the constant sequence  $(x, x, \dots)$ . Whereas in (a) this is the limit of every point in  $(0, 1)$ , in the other cases these fixed points are isolated in the sense of Definition 2.3. The isolated points, however, give rise to sequences that converge to more than one point in the form of limit cycles as shown in (b)–(d).

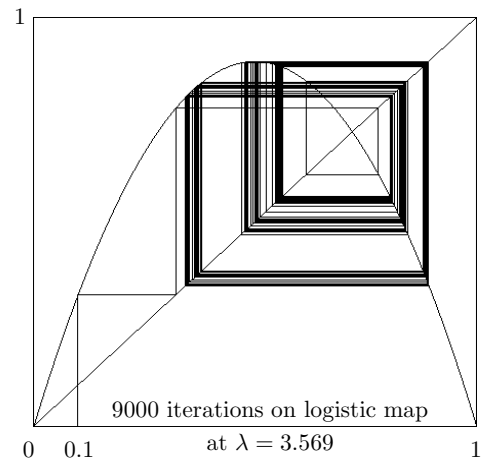
(the multifunctions are generated only in the limiting sense of  $i \rightarrow \infty$  and represent a boundary between functional and non-functional relations on a set), decrease in magnitude with increasing iterations until they reduce to points. This gives rise to a (totally disconnected) Cantor set on the  $y$ -axis in contrast with the connected intervals that the multifunctional limits at  $\lambda > \lambda_*$  of Figs. 16(b)–16(d) produce. By our characterization Definition 4.1 of

chaos therefore,  $\lambda x(1 - x)$  is chaotic for the values of  $\lambda > \lambda_*$  that are shown in Fig. 16. We return to this case in the following subsection.

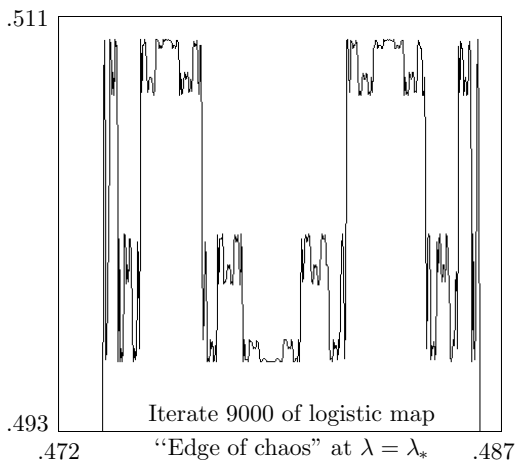
As an example of chaos *in a noniterative system*, we investigate the following question: While maximality of non-injectiveness produced by an increasing number of injective branches is necessary for a family of functions to be chaotic, is this also sufficient for the system to be chaotic? This is an



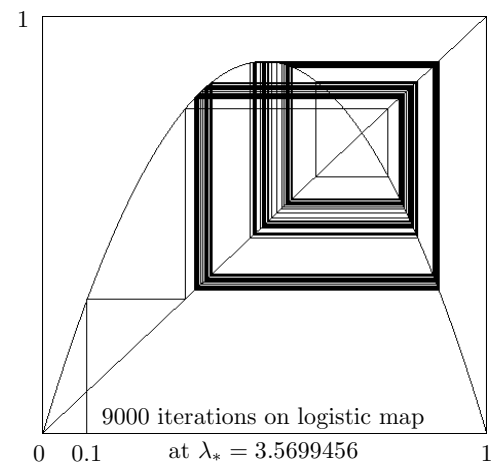
(a)



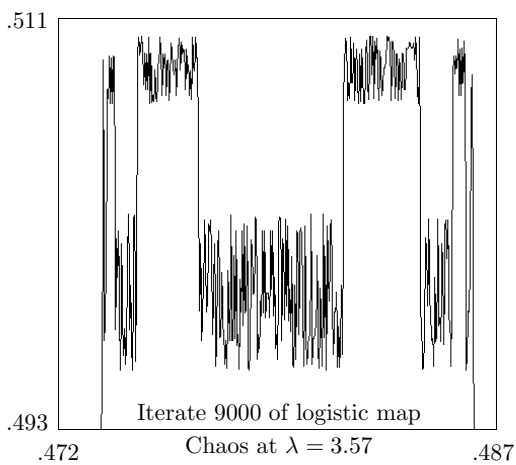
(b)



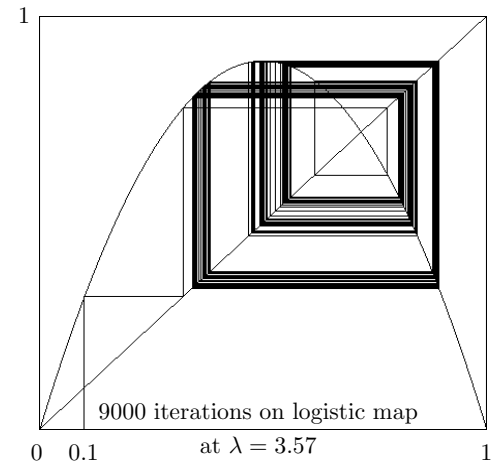
(c)



(d)



(e)



(f)

Fig. 13. Multifunctional and cobweb plots of  $\lambda x(1 - x)$ . Comparison of the graphs for the three values of  $\lambda$  shown in (a)–(f) illustrates how the dramatic changes in the character of the former are conspicuously absent in the conventional plots that display no perceptible distinction between the three cases.

important question especially in the context of a non-iterative family of functions where fixed points are no longer relevant.

Consider the sequence of functions  $|\sin(\pi n x)|_{n=1}^\infty$ . The graphs of the subsequence  $|\sin(2^{n-1}\pi x)|$  and of the sequence  $(t^n(x))$  on  $[0, 1]$  are qualitatively similar in that they both contain  $2^{n-1}$  of their functional graphs each on a base of  $1/2^{n-1}$ . Thus both  $|\sin(2^{n-1}\pi x)|_{n=1}^\infty$  and  $(t^n(x))_{n=1}^\infty$  converge graphically to the multifunction  $[0,1]$  on the same set of points equivalent to 0. This is sufficient for us to conclude that  $|\sin(2^{n-1}\pi x)|_{n=1}^\infty$ , and hence  $|\sin(\pi n x)|_{n=1}^\infty$ , is chaotic on the infinite equivalent set  $[0]$ . While Fig. 9 was a comparison of the first four iterates of the tent and absolute sine maps, Fig. 14 shows the “converged” graphical limits after 17 iterations.

### 4.1. The chaotic attractor

One of the most fascinating characteristics of chaos in dynamical systems is the appearance of attractors the dynamics on which are chaotic. For a subset  $A$  of a topological space  $(X, \mathcal{U})$  such that  $\mathcal{R}(f(A))$  is contained in  $A$  — in this section, unless otherwise stated to the contrary,  $f(A)$  will denote the graph and not the range (image) of  $f$  — which ensures that the iteration process can be carried out in  $A$ , let

$$\begin{aligned} f_{\mathbb{R}_i}(A) &= \bigcup_{j \geq i \in \mathbb{N}} f^j(A) \\ &= \bigcup_{j \geq i \in \mathbb{N}} \left( \bigcup_{x \in A} f^j(x) \right) \end{aligned} \tag{57}$$

generate the filter-base  ${}_F\mathcal{B}$  with  $A_i := f_{\mathbb{R}_i}(A) \in {}_F\mathcal{B}$  being decreasingly nested,  $A_{i+1} \subseteq A_i$  for all  $i \in \mathbb{N}$ , in accordance with Definition A.1.1. The existence of a maximal chain with a corresponding maximal element as assured by the Hausdorff Maximal Principle and Zorn’s Lemma respectively implies a nonempty core of  ${}_F\mathcal{B}$ . As in Sec. 3 following Definition 3.3, we now identify the filterbase with the neighborhood base at  $f^\infty$  which allows us to define

$$\begin{aligned} \text{Atr}(A_1) &\stackrel{\text{def}}{=} \text{adh}({}_F\mathcal{B}) \\ &= \bigcap_{A_i \in {}_F\mathcal{B}} \text{Cl}(A_i) \end{aligned} \tag{58}$$

as the attractor of the set  $A_1$ , where the last equality follows from Eqs. (59) and (20) and the closure is with respect to the topology induced by the neighborhood filter base  ${}_F\mathcal{B}$ . Clearly the attractor as defined here is the graphical limit of the sequence of

functions  $(f^i)_{i \in \mathbb{N}}$  which may be verified by reference to Definition A.1.8, Theorem A.1.3 and the proofs of Theorems A.1.4 and A.1.5, together with the directed set Eq. (A.10) with direction (A.11). The *basin of attraction* of the attractor is  $A_1$  because the graphical limit  $(\mathcal{D}_+, F(\mathcal{D}_+)) \cup (G(\mathcal{R}_+), \mathcal{R}_+)$  of Definition 3.1 may be obtained, as indicated above, by a proper choice of sequences associated with  $\mathcal{A}$ . Note that in the context of iterations of functions, the graphical limit  $(\mathcal{D}_+, y_0)$  of the sequence  $(f^n(x))$  denotes a stable fixed point  $x_*$  with image  $x_* = f(x_*) = y_0$  to which iterations starting at any point  $x \in \mathcal{D}_+$  converge. The graphical limits  $(x_{i0}, \mathcal{R}_+)$  are generated with respect to the class  $\{x_{i*}\}$  of points satisfying  $f(x_{i0}) = x_{i*}$ ,  $i = 0, 1, 2, \dots$  equivalent to unstable fixed point  $x_* := x_{0*}$  to which inverse iterations starting at any initial point in  $\mathcal{R}_+$  must converge. Even though only  $x_*$  is inverse stable, an equivalent class of graphically converged limit multis is produced at every member of the class  $x_{i*} \in [x_*]$ , resulting in the far-reaching consequence that every member of the class is as significant as the parent fixed point  $x_*$  from which they were born in determining the dynamics of the evolving system. The point to remember about infinite intersections of a collection of sets having finite intersection property, as in Eq. (58), is that this may very well be empty; recall, however, that in a compact space this is guaranteed not to be so. In the general case, if  $\text{core}(\mathcal{A}) \neq \emptyset$  then  $\mathcal{A}$  is the principal filter at this core, and  $\text{Atr}(A_1)$  by Eqs. (58) and (A.33) is the closure of this core, which in this case of topology being induced by the filterbase, is just the core itself.  $A_1$  by its very definition, is a positively invariant set as any sequence of graphs converging to  $\text{Atr}(A_1)$  must be eventually in  $A_1$ : the entire sequence therefore lies in  $A_1$ . Clearly, from Theorem A.3.1 and its corollary, the attractor is a positively invariant compact set. A typical attractor is illustrated by the derived sets in the second column of Fig. 22 which also illustrates that the set of functional relations are open in  $\text{Multi}(X)$ ; specifically functional–non-functional correspondences are neutral-selfish related as in Fig. 22, 3–2, with the attracting graphical limit of Eq. (58) forming the boundary of (finitely) many-to-one functions and the one-to-(finitely) many multifunctions.

Equation (58) is to be compared with the *image definition of an attractor* [Stuart & Humphries, 1996] where  $f(A)$  denotes the range and not the graph of  $f$ . Then Eq. (58) can be used to define a

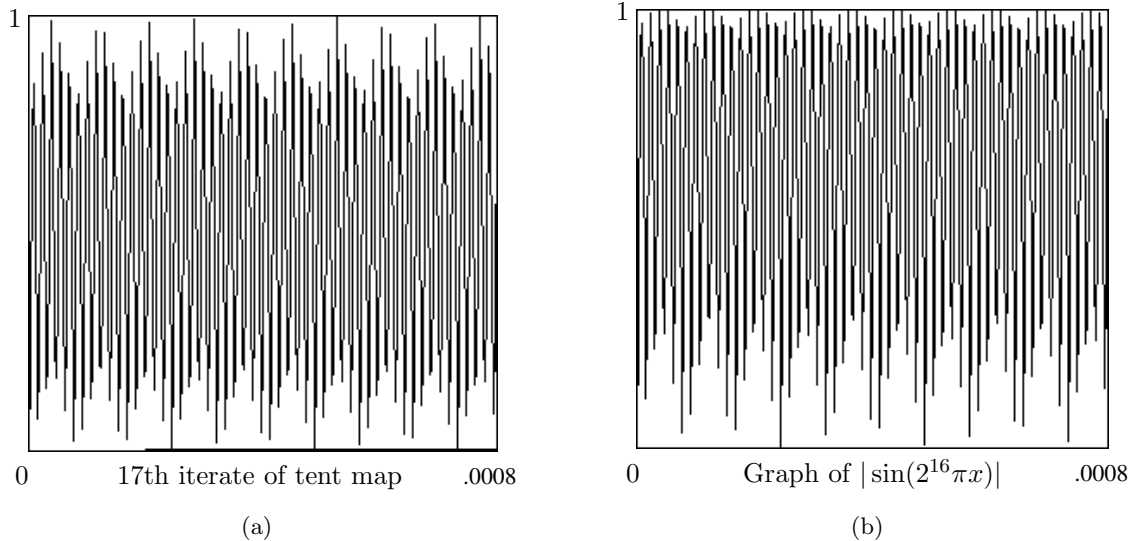


Fig. 14. Similarity in the behavior of the graphs of (a) tent and (b)  $|\sin(2^{16}\pi x)|$  maps at 17 iterations demonstrate chaoticity of the latter.

sequence of points  $x_k \in A_{n_k}$  and hence the subset

$$\begin{aligned} \omega(A) &\stackrel{\text{def}}{=} \{x \in X : (\exists n_k \in \mathbb{N})(n_k \rightarrow \infty)(\exists x_k \in A_{n_k}) \\ &\quad (f^{n_k}(x_k) \rightarrow x)\} \\ &= \{x \in X : (\forall N \in \mathcal{N}_x)(\forall A_i \in \mathcal{A}) \\ &\quad (N \cap A_i \neq \emptyset)\} \end{aligned} \tag{59}$$

as the corresponding attractor of  $A$  that satisfies an equation formally similar to (58) with the difference that the filter-base  $\mathcal{A}$  is now in terms of the image  $f(A)$  of  $A$ , which allows the adherence expression to take the particularly simple form

$$\omega(A) = \bigcap_{i \in \mathbb{N}} \text{Cl}(f^i(A)). \tag{60}$$

The complimentary subset excluded from this definition of  $\omega(A)$ , as compared to  $\text{Atr}(A_1)$ , that is required to complete the formalism is given by Eq. (61) below. Observe that the equation for  $\omega(A)$  is essentially Eq. (A.15), even though we prefer to use the alternate form of Eq. (A.16) as this brings out more clearly the frequenting nature of the sequence. The basin of attraction

$$\begin{aligned} B_f(A) &= \{x \in A : \omega(x) \subseteq \text{Atr}(A)\} \\ &= \{x \in A : (\exists n_k \in \mathbb{N})(n_k \rightarrow \infty) \\ &\quad (f^{n_k}(x) \rightarrow x^* \in \omega(A))\} \end{aligned} \tag{61}$$

of the attractor is the smallest subset of  $X$  in which sequences generated by  $f$  must eventually lie in order to adhere at  $\omega(A)$ . Comparison of Eqs. (62) with (33) and (61) with (32) show that  $\omega(A)$  can

be identified with the subset  $\mathcal{R}_+$  on the  $y$ -axis on which the multifunctional limits  $G : \mathcal{R}_+ \rightarrow X$  of graphical convergence are generated, with its basin of attraction being contained in the  $\mathcal{D}_+$  associated with the injective branch of  $f$  that generates  $\mathcal{R}_+$ . In summary it may be concluded that since definitions (59) and (61) involve both the domain and range of  $f$ , a description of the attractor in terms of the graph of  $f$ , like that of Eq. (58), is more pertinent and meaningful as it combines the requirements of both these equations. Thus, for example, as  $\omega(A)$  is not the function  $G(\mathcal{R}_+)$ , this attractor does not include the equivalence class of inverse stable points that may be associated with  $x_*$ , see for example Fig. 15.

From Eq. (59), we may make the particularly simple choice of  $(x_k)$  to satisfy  $f^{n_k}(x_{-k}) = x$  so that  $x_{-k} = f_B^{-n_k}(x)$ , where  $x_{-k} \in [x_{-k}] := f^{-n_k}(x)$  is the element of the equivalence class of the inverse image of  $x$  corresponding to the injective branch  $f_B$ . This choice is of special interest to us as it is the class that generates the  $G$ -function on  $\mathcal{R}_+$  in graphical convergence. This allows us to express  $\omega(A)$  as

$$\begin{aligned} \omega(A) &= \{x \in X : (\exists n_k \in \mathbb{N})(n_k \rightarrow \infty)(f_B^{-n_k}(x) \\ &\quad = x_{-k} \text{ converges in } (X, \mathcal{U}))\}; \end{aligned} \tag{62}$$

note that the  $x_{-k}$  of this equation and the  $x_k$  of Eq. (59) are, in general, quite different points.

A simple illustrative example of the construction of  $\omega(A)$  for the positive injective branch of the homeomorphism  $(4x^2 - 1)/3$ ,  $-1 \leq x \leq 1$ , is

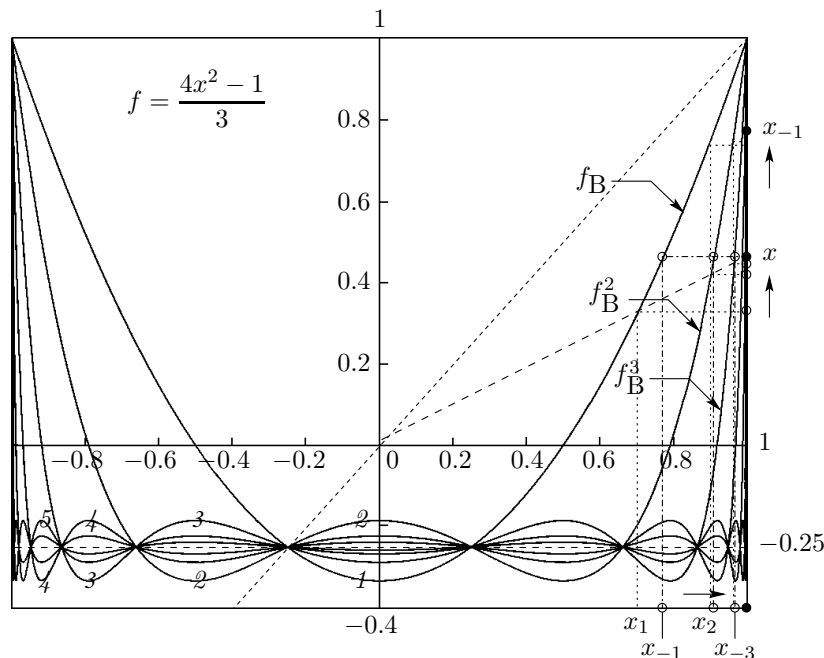


Fig. 15. The attractor for  $f(x) = (4x^2 - 1)/3$ , for  $-1 \leq x \leq 1$ . The converging sequences are denoted by arrows on the right, and  $(x_k)$  are chosen according to the construction shown. This example demonstrates how although  $A \subseteq f(A)$ , where  $A = [0, 1]$  is the domain of the positive injective branch of  $f$ , the succeeding images  $(f^i(A))_{i \geq 1}$  satisfy the required restriction for iteration, and  $A$  in the discussion above can be taken to be  $f(A)$ ; this is permitted as only a finite number of iterates is thereby discarded. It is straightforward to verify that  $\text{Atr}(A_1) = (-1, [-0.25, 1]) \cup ((-1, 1), -0.25) \cup (1, [-0.25, 1])$  with  $F(x) = -0.25$  on  $\mathcal{D}_- = (-1, 1) = \mathcal{D}_+$  and  $G(y) = 1$ , and  $-1$  on  $\mathcal{R}_- = [-0.25, 1] = \mathcal{R}_+$ . By comparison,  $\omega(A)$  from either its definition Eq. (59) or from the equivalent intersection expression Eq. (60), is simply the closed interval  $\mathcal{R}_+ = [-0.25, 1]$ . The italicized iterate numbers on the graphs show how the oscillations die out with increasing iterations from  $x = \pm 1$  and approach  $-0.25$  in all neighborhoods of 0.

shown in Fig. 15, where the arrow-heads denote the converging sequences  $f^{n_i}(x_i) \rightarrow x$  and  $f^{n_i-m}(x_i) \rightarrow x_{-m}$  which proves invariance of  $\omega(A)$  for a homeomorphic  $f$ ; here continuity of the function and its inverse is explicitly required for invariance. Positive invariance of a subset  $A$  of  $X$  implies that for any  $n \in \mathbb{N}$  and  $x \in A$ ,  $f^n(x) = y_n \in A$ , while negative invariance assures that for any  $y \in A$ ,  $f^{-n}(y) = x_{-n} \in A$ . Invariance of  $A$  in both the forward and backward directions therefore means that for any  $y \in A$  and  $n \in \mathbb{N}$ , there exists a  $x \in A$  such that  $f^n(x) = y$ . In interpreting this figure, it may be useful to recall from Definition 4.1 that an increasing number of injective branches of  $f$  is a necessary, but not sufficient, condition for the occurrence of chaos; thus in Figs. 12(a) and 15, increasing noninjectivity of  $f$  leads to constant valued limit functions over a connected  $\mathcal{D}_+$  in a manner similar to that associated with the classical Gibb’s phenomenon in the theory of Fourier series.

Graphical convergence of an increasingly nonlinear family of functions implied by its increasing non-injectivity may now be combined with the re-

quirements of an attractor to lead to the concept of a chaotic attractor to be that on which the dynamics is chaotic in the sense of Definitions 4.1. and 4.2. Hence

**Definition 4.3** (Chaotic Attractor). Let  $A$  be a positively invariant subset of  $X$ . The attractor  $\text{Atr}(A)$  is chaotic on  $A$  if there is sensitive dependence on initial conditions for all  $x \in A$ . The sensitive dependence manifests itself as multifunctional graphical limits for all  $x \in \mathcal{D}_+$  and as chaotic orbits when  $x \notin \mathcal{D}_+$ .

The picture of chaotic attractors that emerge from the foregoing discussions and our characterization of chaos of Definition 4.1 is that it is a subset of  $X$  that is simultaneously “spiked” multifunctional on the  $y$ -axis and consists of a dense collection of singleton domains of attraction on the  $x$ -axis. This is illustrated in Fig. 16 which shows some typical chaotic attractors. The first four diagrams (a)–(d) are for the logistic map with (b)–(d) showing the 4-, 2- and 1-piece attractors for  $\lambda = 3.575$ ,

3.66, and 3.8, respectively that are in qualitative agreement with the standard bifurcation diagram reproduced in (e). Figures 16(b)–16(d) have the advantage of clearly demonstrating how the attractors are formed by considering the graphically converged limit as the object of study unlike in (e) which shows the values of the 501–1001th iterates of  $x_0 = 1/2$  as a function of  $\lambda$ . The difference in (a) and (b) for a change of  $\lambda$  from  $\lambda > \lambda_* = 3.5699456$  to 3.575 is

significant as  $\lambda = \lambda_*$  marks the boundary between the nonchaotic region for  $\lambda < \lambda_*$  and the chaotic for  $\lambda > \lambda_*$  (this is to be understood as being suitably modified by the appearance of the nonchaotic windows for some specific intervals in  $\lambda > \lambda_*$ ). At  $\lambda_*$  the generated fractal Cantor set  $\Lambda$  is an attractor as it attracts almost every initial point  $x_0$  so that the successive images  $x^n = f^n(x_0)$  converge toward the Cantor set  $\Lambda$ . In (f) the chaotic attractors for

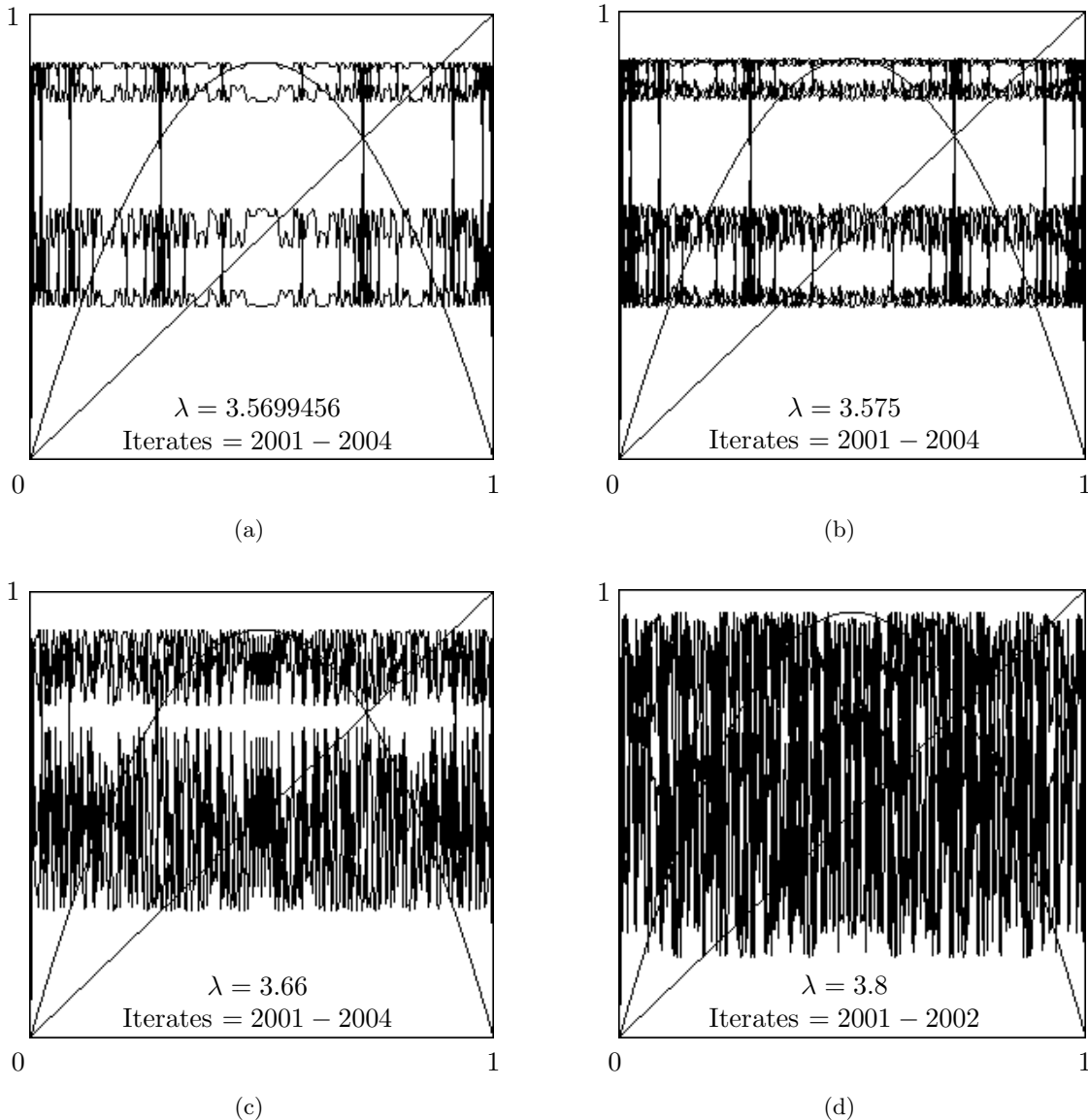


Fig. 16. Chaotic attractors for different values of  $\lambda$ . For the logistic map the usual bifurcation (e) shows the chaotic attractors for  $\lambda > \lambda_* = 3.5699456$ , while (a)–(d) display the graphical limits for four values of  $\lambda$  chosen for the Cantor set and 4-, 2-, and 1-piece attractors, respectively. In (f) the attractor  $[0, 1]$  (where the dotted lines represent odd iterates and the solid lines even iterates of  $f$ ) disappears if  $f$  is reflected about the  $x$ -axis. The function  $f_f(x)$  is given by

$$f_f(x) = \begin{cases} 2(1+x)/3 & 0 \leq x < 1/2 \\ 2(1-x) & 1/2 \leq x \leq 1. \end{cases}$$

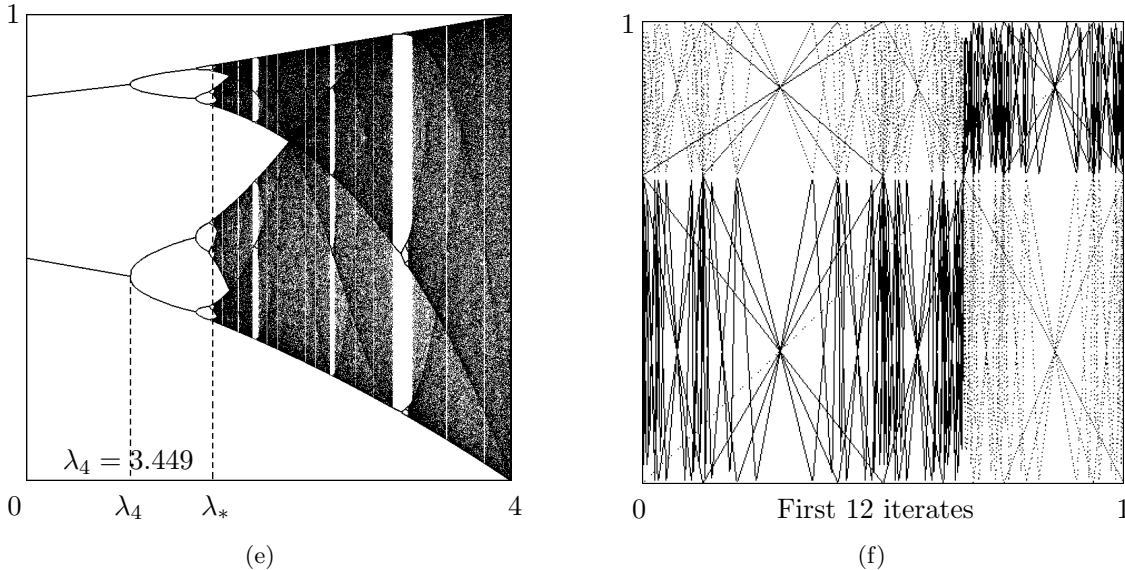


Fig. 16. (Continued)

the piecewise continuous function on  $[0, 1]$

$$f_t(x) = \begin{cases} \frac{2(1+x)}{3}, & 0 \leq x < \frac{1}{2} \\ 2(1-x), & \frac{1}{2} \leq x \leq 1, \end{cases}$$

is  $[0, 1]$  where the dotted lines represent odd iterates and the full lines even iterates of  $f$ ; here the attractor disappears if the function is reflected about the  $x$ -axis.

#### 4.2. Why chaos? A preliminary inquiry

The question as to why a natural system should evolve chaotically is both interesting and relevant, and this section attempts to advance a plausible answer to this inquiry that is based on the connection between topology and convergence contained in the Corollary to Theorem A.1.5. Open sets are groupings of elements that govern convergence of nets and filters, because the required property of being either eventually or frequently in (open) neighborhoods of a point determines the eventual behavior of the net; recall in this connection the unusual convergence characteristics in cofinite and cocountable spaces. Conversely for a given convergence characteristic of a class of nets, it is possible to infer the topology of the space that is responsible for this convergence, and it is this point of view that we adopt here to investigate the question of this subsection: recall that our Definitions 4.1 and 4.2 were based on purely algebraic set-theoretic arguments

on ordered sets, just as the role of the choice of an appropriate problem-dependent basis was highlighted at the end of Sec. 2. Chaos as manifest in its attractors is a direct consequence of the increasing nonlinearity of the map with increasing iteration; we reemphasize that this is only a necessary condition so that the increasing nonlinearities of Figs. 12 and 15 eventually lead to stable states and not to chaotic instability. Under the right conditions as enunciated following Fig. 10, chaos appears to be the natural outcome of the difference in the behavior of a function  $f$  and its inverse  $f^{-}$  under their successive applications. Thus  $f = ff^{-}f$  allows  $f$  to take advantage of its multi-inverse to generate all possible equivalence classes that are available, a feature not accessible to  $f^{-} = f^{-}ff^{-}$ . As we have seen in the foregoing, equivalence classes of fixed points, stable and unstable, are of defining significance in determining the ultimate behavior of an evolving dynamical system and as the eventual (as also frequent) character of a filter or net in a set is dictated by open neighborhoods of points of the set, *it is postulated that chaoticity on a set  $X$  leads to a reformulation of the open sets of  $X$  to equivalence classes generated by the evolving map  $f$* , see Example 2.4(3). Such a redefinition of open sets of equivalence classes allow the evolving system to temporally access an ever increasing number of states even though the equivalent fixed points are not fixed under iterations of  $f$  except for the parent of the class, and can be considered to be the governing criterion for the cooperative or collective behavior



of the system. The predominance of the role of  $f^-$  in  $f = ff^-f$  in generating the equivalence classes (that is exploiting the many-to-one character) of  $f$  is reflected as limit multis for  $f$  (i.e. constant  $f^-$  on  $\mathcal{R}_+$ ) in  $f^- = f^-ff^-$ ; this interpretation of the dynamics of chaos is meaningful as graphical convergence leading to chaos is a result of pointwise bi-convergence of the sequence of iterates of the functions generated by  $f$ . But as  $f$  is a noninjective function on  $X$  possessing the property of increasing nonlinearity in the form of increasing noninjectivity with iteration, various cycles of disjoint equivalence classes are generated under iteration, see for example Fig. 9(a) for the tent map. A reference to Fig. shows that the basic set  $X_B$ , for a finite number  $n$  of iterations of  $f$ , contains the parent of each of these open equivalent sets in the domain of  $f$ , with the topology on  $X_B$  being the corresponding  $p$ -images of these disjoint saturated open sets of the domain. In the limit of infinite iterations of  $f$  leading to the multifunction  $\mathcal{M}$  (this is the  $f^\infty$  of Sec. 4.1), the generated open sets constitute a basis for a topology on  $\mathcal{D}(f)$  and the basis for the topology of  $\mathcal{R}(f)$  are the corresponding  $\mathcal{M}$ -images of these equivalent classes. *It is our contention that the motive force behind evolution toward a chaos, as defined by Definition 4.1, is the drive toward a state of the dynamical system that supports ininality of the limit multi  $\mathcal{M}$* ; see Appendix A.2 with the discussions on Fig. and Eq. (26) in Sec. 2. In the limit of infinite iterations therefore, the open sets of the range  $\mathcal{R}(f) \subseteq X$  are the multi images that graphical convergence generates at each of these inverse-stable fixed points.  $X$  therefore has two topologies imposed on it by the dynamics of  $f$ : the first of equivalence classes generated by the limit multi  $\mathcal{M}$  in the domain of  $f$  and the second as  $\mathcal{M}$ -images of these classes in the range of  $f$ . Quite clearly these two topologies need not be the same; their intersection therefore can be defined to be the *chaotic topology on  $X$  associated with the chaotic map  $f$  on  $X$* . Neighborhoods of points in this topology cannot be arbitrarily small as they consist of all members of the equivalence class to which any element belongs; hence a sequence converging to any of these elements necessarily converges to all of them, and the eventual objective of chaotic dynamics is to generate a topology in  $X$  with respect to which elements of the set can be grouped together in as large equivalence classes as possible in the sense that if a net converges simultaneously

to points  $x \neq y \in X$  then  $x \sim y$ :  $x$  is of course equivalent to itself while  $x, y, z$  are equivalent to each other iff they are simultaneously in every open set in which the net may eventually belong. This hall-mark of chaos can be appreciated in terms of a necessary obliteration of any separation property that the space might have originally possessed, see property (H3) in Appendix A.3. We reemphasize that a set in this chaotic context is required to act in a dual capacity depending on whether it carries the initial or final topology under  $\mathcal{M}$ .

This preliminary inquiry into the nature of chaos is concluded in the final section of this paper.

### 5. Graphical Convergence Works

We present in this section some real evidence in support of our hypothesis of graphical convergence of functions in  $\text{Multi}(X, Y)$ . The example is taken from neutron transport theory, and concerns the discretized spectral approximation [Sengupta, 1988, 1995] of Case's singular eigenfunction solution of the monoenergetic neutron transport equation, [Case & Zweifel, 1967]. The neutron transport equation is a linear form of the Boltzmann equation that is obtained as follows. Consider the neutron-moderator system as a mixture of two species of gases each of which satisfies a Boltzmann equation of the type

$$\begin{aligned} & \left( \frac{\partial}{\partial t} + v_i \cdot \nabla \right) f_i(r, v, t) \\ &= \int dv' \int dv_1 \int dv'_1 \sum_j W_{ij}(v_i \rightarrow v'; v_1 \rightarrow v'_1) \\ & \quad \{ f_i(r, v', t) f_j(r, v'_1, t) - f_i(r, v, t) f_j(r, v_1, t) \} \end{aligned}$$

where

$$W_{ij}(v_i \rightarrow v'; v_1 \rightarrow v'_1) = |v - v_1| \sigma_{ij}(v - v', v_1 - v'_1)$$

$\sigma_{ij}$  being the cross-section of interaction between species  $i$  and  $j$ . Denote neutrons by subscript 1 and the background moderator with which the neutrons interact by 2, and make the assumptions that

- (i) The neutron density  $f_1$  is much less compared with that of the moderator  $f_2$  so that the terms  $f_1 f_1$  and  $f_1 f_2$  may be neglected in the neutron and moderator equations, respectively.
- (ii) The moderator distribution  $f_2$  is not affected by the neutrons. This decouples the neutron and moderator equations and leads to an equilibrium

Maxwellian  $f_M$  for the moderator while the neutrons are described by the linear equation

$$\begin{aligned} & \left( \frac{\partial}{\partial t} + v \cdot \nabla \right) f(r, v, t) \\ &= \int dv' \int dv_1 \int dv'_1 W_{12}(v \rightarrow v'; v_1 \rightarrow v'_1) \\ & \quad \{f(r, v', t)f_M(v'_1) - f(r, v, t)f_M(v_1)\} \end{aligned}$$

This is now put in the standard form of the neutron transport equation [Williams, 1971]

$$\begin{aligned} & \left( \frac{1}{v} \frac{\partial}{\partial t} + \Omega \cdot v + \mathcal{S}(E) \right) \Phi(r, E, \hat{\Omega}, t) \\ &= \int d\Omega' \int dE' \mathcal{S}(r, E' \rightarrow E; \hat{\Omega}' \cdot \hat{\Omega}) \Phi(r, E', \hat{\Omega}', t). \end{aligned}$$

where  $E = mv^2/2$  is the energy and  $\hat{\Omega}$  the direction of motion of the neutrons. The steady state, monoenergetic form of this equation is Eq. (A.53)

$$\begin{aligned} \mu \frac{\partial \Phi(x, \mu)}{\partial x} + \Phi(x, \mu) &= \frac{c}{2} \int_{-1}^1 \Phi(x, \mu') d\mu', \\ 0 < c < 1, \quad -1 \leq \mu \leq 1 \end{aligned}$$

and its singular eigenfunction solution for  $x \in (-\infty, \infty)$  is given by Eq. (A.56)

$$\begin{aligned} \Phi(x, \mu) &= a(\nu_0) e^{-x/\nu_0} \phi(\mu, \nu_0) \\ &+ a(-\nu_0) e^{x/\nu_0} \phi(-\nu_0, \mu) \\ &+ \int_{-1}^1 a(\nu) e^{-x/\nu} \phi(\mu, \nu) d\nu; \end{aligned}$$

see Appendix A.4 for an introductory review of Case’s solution of the one-speed neutron transport equation.

The term “eigenfunction” is motivated by the following considerations. Consider the eigenvalue equation

$$(\mu - \nu)\mathcal{F}_\nu(\mu) = 0, \quad \mu \in V(\mu), \quad \nu \in \mathbb{R} \quad (63)$$

in the space of multifunctions  $\text{Multi}(V(\mu), (-\infty, \infty))$ , where  $\mu$  is in either of the intervals  $[-1, 1]$  or  $[0, 1]$  depending on whether the given boundary conditions for Eq. (A.53) is full-range or half range. If we are looking only for functional solutions of Eq. (63), then the unique function  $\mathcal{F}$  that satisfies this equation for all possible  $\mu \in V(\mu)$  and  $\nu \in \mathbb{R} - V(\mu)$  is  $\mathcal{F}_\nu(\mu) = 0$  which means, according to Table 1, that the point spectrum of  $\mu$  is empty and  $(\mu - \nu)^{-1}$  exists for all  $\nu$ . When  $\nu \in V(\mu)$ , however, this inverse is not continuous and we show below that in  $\text{Map}(V(\mu), 0)$ ,  $\nu \in V(\mu)$  belongs to

the continuous spectrum of  $\mu$ . This distinction between the nature of the inverses depending on the relative values of  $\mu$  and  $\nu$  suggests a wider “non-function” space in which to look for the solutions of operator equations, and in keeping with the philosophy embodied in Fig. of treating inverse problems in the space of multifunctions, we consider all  $\mathcal{F}_\nu \in \text{Multi}(V(\mu), \mathbb{R})$  satisfying Eq. (63) to be eigenfunctions of  $\mu$  for the corresponding eigenvalue  $\nu$ , leading to the following multifunctional solution of (63)

$$\mathcal{F}_\nu(\mu) = \begin{cases} (V(\mu), 0) & \text{if } \nu \notin V(\mu) \\ (V(\mu) - \nu, 0) \cup (\nu, \mathbb{R}) & \text{if } \nu \in V(\mu), \end{cases}$$

where  $V(\mu) - \nu$  is used as a shorthand for the interval  $V(\mu)$  with  $\nu$  deleted. Rewriting the eigenvalue equation (63) as  $\mu_\nu(\mathcal{F}_\nu(\mu)) = 0$  and comparing this with Fig. , allows us to draw the correspondences

$$\begin{aligned} f &\Leftrightarrow \mu_\nu \\ X \text{ and } Y &\Leftrightarrow \{ \mathcal{F}_\nu \in \text{Multi}(V(\mu), \mathbb{R}) : \\ & \quad \mathcal{F}_\nu \in \mathcal{D}(\mu_\nu) \} \\ f(X) &\Leftrightarrow \{ 0 : 0 \in Y \} \\ X_B &\Leftrightarrow \{ 0 : 0 \in X \} \\ f^- &\Leftrightarrow \mu_\nu^- . \end{aligned} \quad (64)$$

Thus a multifunction in  $X$  is equivalent to 0 in  $X_B$  under the linear map  $\mu_\nu$ , and we show below that this multifunction is in fact the Dirac delta “function”  $\delta_\nu(\mu)$ , usually written as  $\delta(\mu - \nu)$ . This suggests that in  $\text{Multi}(V(\mu), \mathbb{R})$ , every  $\nu \in V(\mu)$  is in the point spectrum of  $\mu$ , so that discontinuous functions that are pointwise limits of functions in function space can be replaced by graphically converged multifunctions in the space of multifunctions. Completing the equivalence class of 0 in Fig. , gives the multifunctional solution of Eq. (63).

From a comparison of the definition of ill-posedness (Sec. 2) and the spectrum (Table 1), it is clear that  $\mathcal{L}_\lambda(x) = y$  is ill-posed iff

- (1)  $\mathcal{L}_\lambda$  not injective  $\Leftrightarrow \lambda \in P\sigma(\mathcal{L}_\lambda)$ , which corresponds to the first row of Table 1.
- (2)  $\mathcal{L}_\lambda$  not surjective  $\Leftrightarrow$  the values of  $\lambda$  correspond to the second and third columns of Table 1.
- (3)  $\mathcal{L}_\lambda$  is bijective but not open  $\Leftrightarrow \lambda$  is either in  $C\sigma(\mathcal{L}_\lambda)$  or  $R\sigma(\mathcal{L}_\lambda)$  corresponding to the second row of Table 1.

We verify in the three steps below that  $X = L_1[-1, 1]$  of integrable functions,  $\nu \in V(\mu) = [-1, 1]$  belongs to the continuous spectrum of  $\mu$ .

Table 1. Spectrum of linear operator  $\mathcal{L} \in \text{Map}(X)$ . Here  $\mathcal{L}_\lambda := \mathcal{L} - \lambda$  satisfies the equation  $\mathcal{L}_\lambda(x) = 0$ , with the resolvent set  $\rho(\mathcal{L})$  of  $\mathcal{L}$  consisting of all those complex numbers  $\lambda$  for which  $\mathcal{L}_\lambda^{-1}$  exists as a continuous operator with dense domain. Any value of  $\lambda$  for which this is not true is in the spectrum  $\sigma(\mathcal{L})$  of  $\mathcal{L}$ , that is further subdivided into three disjoint components of the point, continuous and residual spectra according to the criteria shown in the table.

$\mathcal{L}_\lambda$	$\mathcal{L}_\lambda^{-1}$	$\mathcal{R}(\mathcal{L}_\lambda)$		
		$\mathcal{R} = X$	$\text{Cl}(\mathcal{R}) = X$	$\text{Cl}(\mathcal{R}) \neq X$
Not injective	...	$P\sigma(\mathcal{L})$	$P\sigma(\mathcal{L})$	$P\sigma(\mathcal{L})$
Injective	Not continuous	$C\sigma(\mathcal{L})$	$C\sigma(\mathcal{L})$	$R\sigma(\mathcal{L})$
	Continuous	$\rho(\mathcal{L})$	$\rho(\mathcal{L})$	$R\sigma(\mathcal{L})$

(a)  $\mathcal{R}(\mu_\nu)$  is dense, but not equal to  $L_1$ . The set of functions  $g(\mu) \in L_1$  such that  $\mu_\nu^{-1}g \in L_1$  cannot be the whole of  $L_1$ . Thus, for example, the piecewise constant function  $g = \text{const} \neq 0$  on  $|\mu - \nu| \leq \delta > 0$  and 0 otherwise is in  $L_1$  but not in  $\mathcal{R}(\mu_\nu)$  as  $\mu_\nu^{-1}g \notin L_1$ . Nevertheless for any  $g \in L_1$ , we may choose the sequence of functions

$$g_n(\mu) = \begin{cases} 0, & \text{if } |\mu - \nu| \leq 1/n \\ g(\mu), & \text{otherwise} \end{cases}$$

in  $\mathcal{R}(\mu_\nu)$  to be eventually in every neighborhood of  $g$  in the sense that  $\lim_{n \rightarrow \infty} \int_{-1}^1 |g - g_n| = 0$ .

(b) The inverse  $(\mu - \nu)^{-1}$  exists but is not continuous. The inverse exists because, as noted earlier, 0 is the only functional solution of Eq. (63).

Nevertheless although the net of functions

$$\delta_{\nu\varepsilon}(\mu) = \frac{1}{\tan^{-1}(1 + \nu)/\varepsilon + \tan^{-1}(1 - \nu)/\varepsilon} \times \left( \frac{\varepsilon}{(\mu - \nu)^2 + \varepsilon^2} \right), \quad \varepsilon > 0$$

is in the domain of  $\mu_\nu$  because  $\int_{-1}^1 \delta_{\nu\varepsilon}(\mu) d\mu = 1$  for all  $\varepsilon > 0$ ,

$$\lim_{\varepsilon \rightarrow 0} \int_{-1}^1 |\mu - \nu| \delta_{\nu\varepsilon}(\mu) d\mu = 0$$

implying that  $(\mu - \nu)^{-1}$  is unbounded.

Taken together, (a) and (b) show that functional solutions of Eq. (63) lead to state 2-2 in Table 1; hence  $\nu \in [-1, 1] = C\sigma(\mu)$ .

(c) The two integral constraints in (b) also mean that  $\nu \in C\sigma(\mu)$  is a *generalized eigenvalue* of  $\mu$  which justifies calling the graphical limit  $\delta_{\nu\varepsilon}(\mu) \xrightarrow{\mathbf{G}} \delta_\nu(\mu)$  a *generalized, or singular, eigen-*

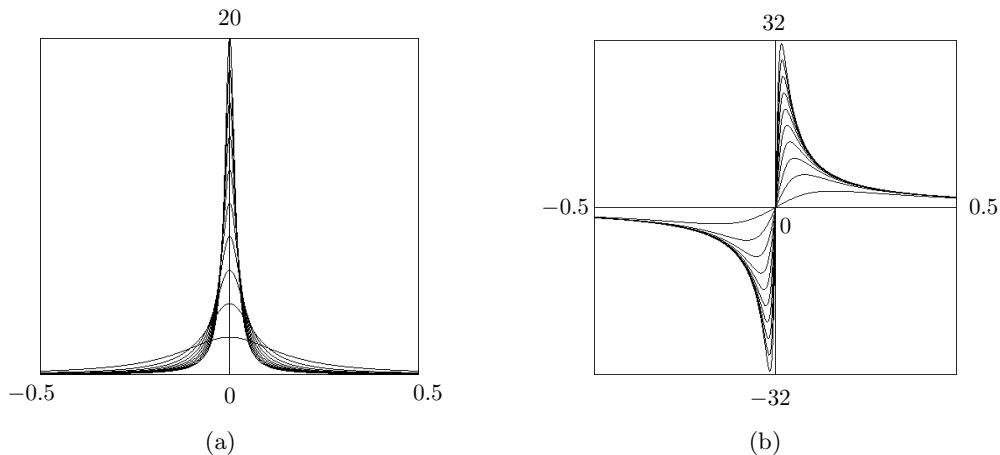


Fig. 17. Graphical convergence of: (a) Poisson kernel  $\delta_\varepsilon(x) = \varepsilon/\pi(x^2 + \varepsilon^2)$  and (b) conjugate Poisson kernel  $P_\varepsilon(x) = x/(x^2 + \varepsilon^2)$  to the Dirac delta and principal value, respectively; the graphs, each for a definite  $\varepsilon$ -value, converges to the respective limits as  $\varepsilon \rightarrow 0$ .

function, see Fig. 17 which clearly indicates the convergence of the net of functions.<sup>26</sup>

From the fact that the solution Eq. (A.56) of the transport equation contains an integral involving the multifunction  $\phi(\mu, \nu)$ , we may draw an interesting physical interpretation. As the multi appears *everywhere* on  $V(\mu)$  (i.e. there are no chaotic orbits but only the multifunctions that produce them), we have here a situation typical of *maximal ill-posedness* characteristic of chaos: note that both the functions comprising  $\phi_\varepsilon(\mu, \nu)$  are non-injective. As the solution (A.56) involves an integral over all  $\nu \in V(\mu)$ , the singular eigenfunctions — that collectively may be regarded as representing a *chaotic substate* of the system represented by the solution of the neutron transport equation — combine with the functional components  $\phi(\pm\nu_0, \mu)$  to produce the well-defined, non-chaotic, experimental end result of the neutron flux  $\Phi(x, \mu)$ .

The solution (A.56) is obtained by assuming  $\Phi(x, \mu) = e^{-x/\nu}\phi(\mu, \nu)$  to get the equation for  $\phi(\mu, \nu)$  to be  $(\mu - \nu)\phi(\mu, \nu) = -c\nu/2$  with the normalization  $\int_{-1}^1 \phi(\mu, \nu) = 1$ . As  $\mu_\nu^{-1}$  is not invertible in  $\text{Multi}(V(\mu), \mathbb{R})$  and  $\mu_{\nu B} : X_B \rightarrow f(X)$  does not exist, the alternate approach of regularization was adopted in [Sengupta, 1988, 1995] to rewrite  $\mu_\nu\phi(\mu, \nu) = -c\nu/2$  as  $\mu_{\nu\varepsilon}\phi_\varepsilon(\mu, \nu) = -c\nu/2$  with  $\mu_{\nu\varepsilon} := \mu - (\nu + i\varepsilon)$  being a net of bijective functions for  $\varepsilon > 0$ ; this is a consequence of the fact that for the multiplication operator every non-real  $\lambda$  belongs to the resolvent set of the operator. The family of solutions of the latter equation is given by [Sengupta 1988, 1995]

$$\phi_\varepsilon(\nu, \mu) = \frac{c\nu}{2} \frac{\nu - \mu}{(\mu - \nu)^2 + \varepsilon^2} + \frac{\lambda_\varepsilon(\nu)}{\pi_\varepsilon} \frac{\varepsilon}{(\mu - \nu)^2 + \varepsilon^2} \quad (65)$$

where the required normalization  $\int_{-1}^1 \phi_\varepsilon(\nu, \mu) = 1$  gives

$$\begin{aligned} \lambda_\varepsilon(\nu) &= \frac{\pi_\varepsilon}{\tan^{-1}(1 + \nu)/\varepsilon + \tan^{-1}(1 - \nu)/\varepsilon} \\ &\times \left( 1 - \frac{c\nu}{4} \ln \frac{(1 + \nu)^2 + \varepsilon^2}{(1 - \nu)^2 + \varepsilon^2} \right) \\ &\xrightarrow{\varepsilon \rightarrow 0} \pi \lambda(\nu) \end{aligned}$$

with

$$\pi_\varepsilon = \varepsilon \int_{-1}^1 \frac{d\mu}{\mu^2 + \varepsilon^2} = 2 \tan^{-1} \left( \frac{1}{\varepsilon} \right) \xrightarrow{\varepsilon \rightarrow 0} \pi.$$

These discretized equations should be compared with the corresponding exact ones of Appendix A.4. We shall see that the net of functions (65) converges graphically to the multifunction Eq. (A.55) as  $\varepsilon \rightarrow 0$ .

In the discretized spectral approximation, the singular eigenfunction  $\phi(\mu, \nu)$  is replaced by  $\phi_\varepsilon(\mu, \nu)$ ,  $\varepsilon \rightarrow 0$ , with the integral in  $\nu$  being replaced by an appropriate sum. The solution Eq. (A.58) of the physically interesting half-space  $x \geq 0$  problem then reduces to [Sengupta, 1988, 1995]

$$\begin{aligned} \Phi_\varepsilon(x, \mu) &= a(\nu_0)e^{-x/\nu_0}\phi(\mu, \nu_0) \\ &+ \sum_{i=1}^N a(\nu_i)e^{-x/\nu_i}\phi_\varepsilon(\mu, \nu_i) \quad \mu \in [0, 1] \end{aligned} \quad (66)$$

where the nodes  $\{\nu_i\}_{i=1}^N$  are chosen suitably. This discretized spectral approximation to Case's solution has given surprisingly accurate numerical results for a set of properly chosen nodes when compared with exact calculations. Because of its involved nature [Case & Zweifel, 1967], the exact calculations are basically numerical which leads to nonlinear integral equations as part of the solution procedure. To appreciate the enormous complexity of the exact treatment of the half-space problem, we recall that the complete set of eigenfunctions  $\{\phi(\mu, \nu_0), \{\phi(\mu, \nu)\}_{\nu \in [0, 1]}\}$  are orthogonal with respect to the half-range weight function  $W(\mu)$  of half-range theory, Eq. (A.61), that is expressed only in terms of solution of the nonlinear integral equation Eq. (A.62). The solution of a half-space problem then evaluates the coefficients  $\{a(\nu_0), a(\nu)_{\nu \in [0, 1]}\}$  from the appropriate half range (that is  $0 \leq \mu \leq 1$ ) orthogonality integrals satisfied by the eigenfunctions  $\{\phi(\mu, \nu_0), \{\phi(\mu, \nu)\}_{\nu \in [0, 1]}\}$  with respect to the weight  $W(\mu)$ , see Appendix A.4 for the necessary details of the half-space problem in neutron transport theory.

As may be appreciated from this brief introduction, solutions to half-space problems are not simple and actual numerical computations must rely a great deal on tabulated values of the  $X$ -function.

<sup>26</sup>The technical definition of a generalized eigenvalue is as follows. Let  $\mathcal{L}$  be a linear operator such that there exists in the domain of  $\mathcal{L}$  a sequence of elements  $(x_n)$  with  $\|x_n\| = 1$  for all  $n$ . If  $\lim_{n \rightarrow \infty} \|(\mathcal{L} - \lambda)x_n\| = 0$  for some  $\lambda \in \mathbb{C}$ , then this  $\lambda$  is a *generalized eigenvalue* of  $\mathcal{L}$ , the corresponding eigenfunction  $x_\infty$  being a *generalized eigenfunction*.

Self-consistent calculations of sample benchmark problems performed by the discretized spectral approximation in a full-range adaption of the half-range problem described below that generate all necessary data, independent of numerical tables, with the quadrature nodes  $\{\nu_i\}_{i=1}^N$  taken at the zero Legendre polynomials show that the full range formulation of this approximation [Sengupta, 1988, 1995] can give very accurate results not only of integrated quantities like the flux  $\Phi$  and leakage of particles out of the half space, but of also basic “raw” data like the extrapolated end point

$$z_0 = \frac{c\nu_0}{4} \int_0^1 \frac{\nu}{N(\nu)} \left(1 + \frac{c\nu^2}{1-\nu^2}\right) \ln \left(\frac{\nu_0 + \nu}{\nu_0 - \nu}\right) d\nu \tag{67}$$

and of the  $X$ -function itself. Given the involved

nature of the exact theory, it is our contention that the remarkable accuracy of these basic data, some of which is reproduced in Table 2, is due to the graphical convergence of the net of functions

$$\phi_\varepsilon(\mu, \nu) \xrightarrow{\mathbf{G}} \phi(\mu, \nu)$$

shown in Fig. 18; here  $\varepsilon = 1/\pi N$  so that  $\varepsilon \rightarrow 0$  as  $N \rightarrow \infty$ . By this convergence, the delta function and principal values in  $[-1, 1]$  are the multifunctions  $([-1, 0), 0) \cup (0, [0, \infty) \cup ((0, 1], 0)$  and  $\{1/x\}_{x \in [-1, 0)} \cup (0, (-\infty, \infty)) \cup \{1/x\}_{x \in (0, 1]}$  respectively. Tables 2 and 3, taken from [Sengupta, 1988] and [Sengupta, 1995], show respectively the extrapolated end point and  $X$ -function by the full-range adaption of the discretized spectral approximation for two different half-range problems denoted as Problems A and B defined as

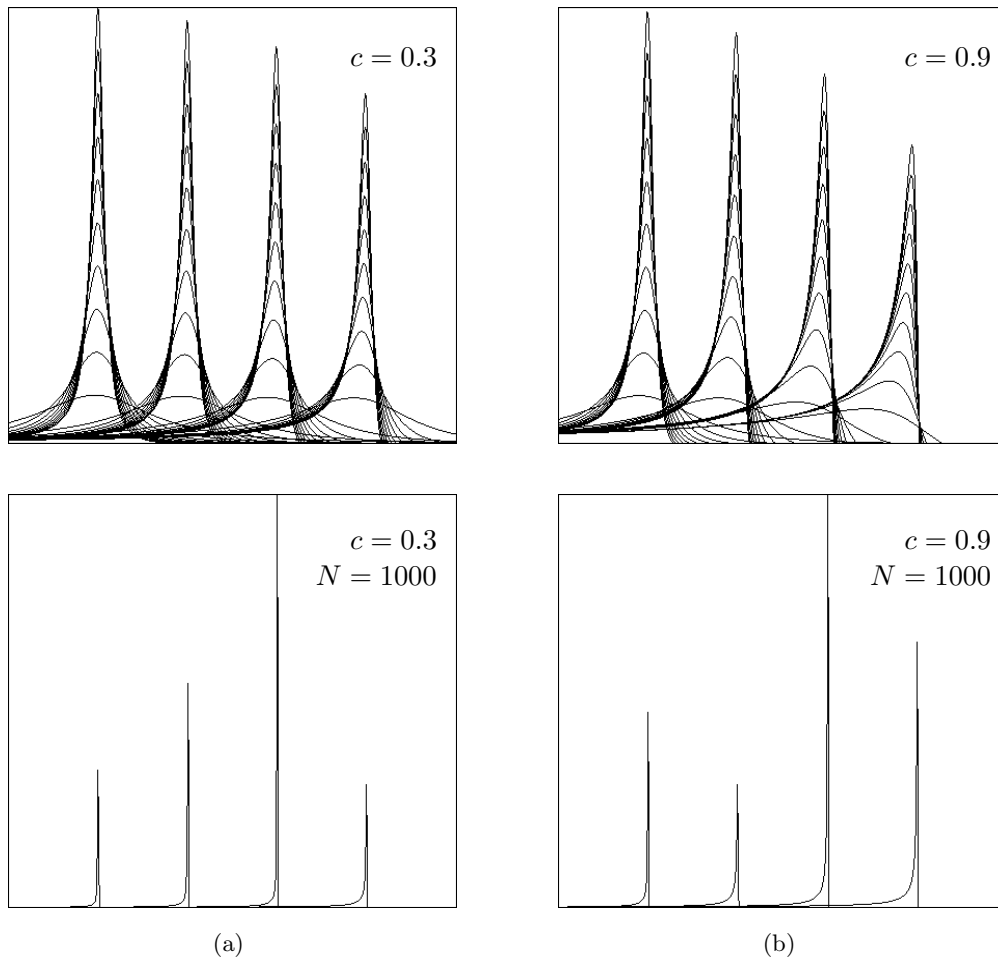


Fig. 18. Rational function approximations  $\phi_\varepsilon(\mu, \nu)$  of the singular eigenfunction  $\phi(\mu, \nu)$  at four different values of  $\nu$ .  $N = 1000$  denotes the “converged” multifunction  $\phi$ , with the peaks at the specific  $\nu$ -values chosen.

- Problem A* Equation :  $\mu\Phi_x + \Phi = (c/2) \int_{-1}^1 \Phi(x, \mu')d\mu'$ ,  $x \geq 0$   
 Boundary condition :  $\Phi(0, \mu) = 0$  for  $\mu \geq 0$   
 Asymptotic condition :  $\Phi \rightarrow e^{-x/\nu_0}\phi(\mu, \nu_0)$  as  $x \rightarrow \infty$ .
- Problem B* Equation :  $\mu\Phi_x + \Phi = (c/2) \int_{-1}^1 \Phi(x, \mu')d\mu'$ ,  $x \geq 0$   
 Boundary condition :  $\Phi(0, \mu) = 1$  for  $\mu \geq 0$   
 Asymptotic condition :  $\Phi \rightarrow 0$  as  $x \rightarrow \infty$ .

The full  $-1 \leq \mu \leq 1$  range form of the half  $0 \leq \mu \leq 1$  range discretized spectral approximation replaces the exact integral boundary condition at  $x = 0$  by a suitable quadrature sum over the values of  $\nu$  taken at the zeros of Legendre polynomials; thus the condition at  $x = 0$  can be expressed as

$$\psi(\mu) = a(\nu_0)\phi(\mu, \nu_0) + \sum_{i=1}^N a(\nu_i)\phi_\varepsilon(\mu, \nu_i), \quad (68)$$

$$\mu \in [0, 1],$$

where  $\psi(\mu) = \Phi(0, \mu)$  is the specified incoming radiation incident on the boundary from the left, and the half-range coefficients  $a(\nu_0), \{a(\nu)\}_{\nu \in [0,1]}$  are to be evaluated using the  $W$ -function of Appendix A.4. We now exploit the relative simplicity of the full-range calculations by replacing Eq. (68) by Eq. (69) following, where the coefficients  $\{b(\nu_i)\}_{i=0}^N$  are used to distinguish the full-range coefficients from the half-range ones. The significance of this change lies in the overwhelming simplicity

of the full-range weight function  $\mu$  as compared to the half-range function  $W(\mu)$ , and the resulting simplicity of the orthogonality relations that follow, see Appendix A.4. The basic data of  $z_0$  and  $X(-\nu)$  are then completely generated self-consistently [Sengupta, 1988, 1995] by the discretized spectral approximation from the full-range adaption

$$\sum_{i=0}^N b_i\phi_\varepsilon(\mu, \nu_i) = \psi_+(\mu) + \psi_-(\mu), \quad (69)$$

$$\mu \in [-1, 1], \nu_i \geq 0$$

Table 2. Extrapolated end-point  $z_0$ .

c	$cz_0$			Exact
	N = 2	N = 6	N = 10	
0.2	0.78478	0.78478	0.78478	0.7851
0.4	0.72996	0.72996	0.72996	0.7305
0.6	0.71535	0.71536	0.71536	0.7155
0.8	0.71124	0.71124	0.71124	0.7113
0.9	0.71060	0.71060	0.71061	0.7106

Table 3.  $X(-\nu)$  by the full-range method.

c	N	$X(-\nu)$			
		$\nu_i$	Problem A	Problem B	Exact
0.2	2	0.2133	0.8873091	0.8873091	0.887308
		0.7887	0.5826001	0.5826001	0.582500
		0.0338	1.3370163	1.3370163	1.337015
		0.1694	1.0999831	1.0999831	1.099983
0.6	6	0.3807	0.8792321	0.8792321	0.879232
		0.6193	0.7215240	0.7215240	0.721524
		0.8306	0.6239109	0.6239109	0.623911
		0.9662	0.5743556	0.5743556	0.574355
		0.0130	1.5971784	1.5971784	1.597163
		0.0674	1.4245314	1.4245314	1.424532
		0.1603	1.2289940	1.2289940	1.228995
0.9	10	0.2833	1.0513750	1.0513750	1.051376
		0.4255	0.9058140	0.9058410	0.905842
		0.5744	0.7934295	0.7934295	0.793430
		0.7167	0.7102823	0.7102823	0.710283
		0.8397	0.6516836	0.6516836	0.651683
		0.9325	0.6136514	0.6136514	0.613653
		0.9870	0.5933988	0.5933988	0.593399

of the discretized boundary condition Eq. (68), where  $\psi_+(\mu)$  is by definition the incident flux  $\psi(\mu)$  for  $\mu \in [0, 1]$  and 0 if  $\mu \in [-1, 0]$ , while

$$\psi_-(\mu) = \begin{cases} \sum_{i=0}^N b_i^- \phi_\varepsilon(\mu, \nu_i) & \text{if } \mu \in [-1, 0], \nu_i \geq 0 \\ 0 & \text{if } \mu \in [0, 1] \end{cases}$$

is the emergent angular distribution out of the medium. Equation (69) corresponds to the full-range  $\mu \in [-1, 1]$ ,  $\nu_i \geq 0$  form

$$\begin{aligned} b(\nu_0)\phi(\mu, \nu_0) + \int_0^1 b(\nu)\phi(\mu, \nu)d\nu \\ = \psi_+(\mu) + \left( b^-(\nu_0)\phi(\mu, \nu_0) + \int_0^1 b^-(\nu)\phi(\mu, \nu)d\nu \right) \end{aligned} \tag{70}$$

of boundary condition (A.59) with the first and second terms on the right having the same interpretation as for Eq. (69). This full-range simulation merely states that the solution (A.58) of Eq. (A.53) holds for all  $\mu \in [-1, 1]$ ,  $x \geq 0$ , although it was obtained, unlike in the regular full-range case, from the given radiation  $\psi(\mu)$  incident on the boundary at  $x = 0$  over only half the interval  $\mu \in [0, 1]$ . To obtain the simulated full-range coefficients  $\{b_i\}$  and  $\{b_i^-\}$  of the half-range problem, we observe that there are effectively only half the number of coefficients as compared to a normal full-range problem because  $\nu$  is now only over half the full interval. This allows us to generate two sets of equations from (70) by integrating with respect to  $\mu \in [-1, 1]$  with  $\nu$  in the half intervals  $[-1, 0]$  and  $[0, 1]$  to obtain the two sets of coefficients  $b^-$  and  $b$ , respectively. Accordingly we get from Eq. (69) with  $j = 0, 1, \dots, N$  the sets of equations

$$\begin{aligned} (\psi, \phi_{j-})_\mu^{(+)} &= - \sum_{i=0}^N b_i^- (\phi_{i+}, \phi_{j-})_\mu^{(-)} \\ b_j &= \frac{1}{N_j} \left( (\psi, \phi_{j+})_\mu^{(+)} + \sum_{i=0}^N b_i^- (\phi_{i+}, \phi_{j+})_\mu^{(-)} \right) \end{aligned} \tag{71}$$

where  $(\phi_{j\pm})_{j=1}^N$  represents  $(\phi_\varepsilon(\mu, \pm\nu_j))_{j=1}^N$ ,  $\phi_{0\pm} = \phi(\mu, \pm\nu_0)$ , the (+) (−) superscripts are used to denote the integrations with respect to  $\mu \in [0, 1]$  and  $\mu \in [-1, 0]$  respectively, and  $(f, g)_\mu$  denotes the usual inner product in  $[-1, 1]$  with respect to the full range weight  $\mu$ . While the first set of  $N + 1$  equations give  $b_i^-$ , the second set produces

the required  $b_j$  from these “negative” coefficients. By equating these calculated  $b_i$  with the exact half-range expressions for  $a(\nu)$  with respect to  $W(\mu)$  as outlined in Appendix A.4, it is possible to find numerical values of  $z_0$  and  $X(-\nu)$ . Thus from the second of Eq. (A.64),  $\{X(-\nu_i)\}_{i=1}^N$  is obtained with  $b_{iB} = a_{iB}$ ,  $i = 1, \dots, N$ , which is then substituted in the second of Eq. (A.63) with  $X(-\nu_0)$  obtained from  $a_A(\nu_0)$  according to Appendix A.4, to compare the respective  $a_{iA}$  with the calculated  $b_{iA}$  from (71). Finally the full-range coefficients of Problem A can be used to obtain the  $X(-\nu)$  values from the second of Eqs. (A.63) and compared with the exact tabulated values as in Table 3. The tabulated values of  $cz_0$  from Eq. (67) show a consistent deviation from our calculations of Problem A according to  $a_A(\nu_0) = -\exp(-2z_0/\nu_0)$ . Since the  $X(-\nu)$  values of Problem A in Table 3 also need the same  $b_{0A}$  as input that was used in obtaining  $z_0$ , it is reasonable to conclude that the “exact” numerical integration of  $z_0$  is inaccurate to the extent displayed in Table 2.

From these numerical experiments and Fig. 18 we may conclude that the continuous spectrum  $[-1, 1]$  of the position operator  $\mu$  acts as the  $\mathcal{D}_+$  points in generating the multifunctional Case singular eigenfunction  $\phi(\mu, \nu)$ . Its rational approximation  $\phi_\varepsilon(\mu, \nu)$  in the context of the simple simulated full-range computations of the complex half-range exact theory of Appendix A.4, clearly demonstrates the utility of graphical convergence of sequence of functions to multifunction. The totality of the multifunctions  $\phi(\mu, \nu)$  for all  $\nu$  in Figs. 18(c) and 18(d) endows the problem with the character of maximal ill-posedness that is characteristic of chaos. This chaotic signature of the transport equation is however latent as the experimental output  $\Phi(x, \mu)$  is well-behaved and regular. This important example shows how nature can use hidden and complex chaotic substates to generate order through a process of superposition.

## 6. Does Nature Support Complexity?

The question of this section is basic in the light of the theory of chaos presented above as it may be reformulated to the inquiry of what makes nature support chaoticity in the form of increasing non-injectivity of an input–output system. It is the purpose of this section to exploit the connection between spectral theory and the dynamics of chaos that has been presented in the previous section.

Since linear operators on finite dimensional spaces do not possess continuous or residual spectra, spectral theory on infinite dimensional spaces essentially involves limiting behavior to infinite dimensions of the familiar matrix eigenvalue–eigenvector problem. As always this means extensions, dense embeddings and completions of the finite dimensional problem that show up as generalized eigenvalues and eigenvectors. In its usual form, the goal of nonlinear spectral theory consists [Appell *et al.*, 2000] in the study of  $T_\lambda^{-1}$  for nonlinear operators  $T_\lambda$  that satisfy more general continuity conditions, like differentiability and Lipschitz continuity, than simple boundedness that is sufficient for linear operators. The following generalization of the concept of the spectrum of a linear operator to the nonlinear case is suggestive. For a nonlinear map,  $\lambda$  need not appear only in a multiplying role, so that an eigenvalue equation can be written more generally as a fixed-point equation

$$f(\lambda; x) = x$$

with a fixed point corresponding to the eigenfunction of a linear operator and an “eigenvalue” being the value of  $\lambda$  for which this fixed point appears. The correspondence of the residual and continuous parts of the spectrum are, however, less trivial than for the point spectrum. This is seen from the following two examples [Roman, 1975]. Let  $Ae_k = \lambda_k e_k, k = 1, 2, \dots$  be an eigenvalue equation with  $e_j$  being the  $j$ th unit vector. Then  $(A - \lambda)e_k := (\lambda_k - \lambda)e_k = 0$  iff  $\lambda = \lambda_k$  so that  $\{\lambda_k\}_{k=1}^\infty \in P\sigma(A)$  are the only eigenvalues of  $A$ . Consider now  $(\lambda_k)_{k=1}^\infty$  to be a sequence of real numbers that tends to a finite  $\lambda^*$ ; for example, let  $A$  be a diagonal matrix having  $1/k$  as its diagonal entries. Then  $\lambda^*$  belongs to the continuous spectrum of  $A$  because  $(A - \lambda^*)e_k = (\lambda_k - \lambda^*)e_k$  with  $\lambda_k \rightarrow \lambda^*$  implies that  $(A - \lambda^*)^{-1}$  is an unbounded linear operator and  $\lambda^*$  a generalized eigenvalue of  $A$ . In the second example  $Ae_k = e_{k+1}/(k + 1)$ , it is not difficult to verify that: (a) The point spectrum of  $A$  is empty, (b) The range of  $A$  is not dense because it does not contain  $e_1$ , and (c)  $A^{-1}$  is unbounded because  $Ae_k \rightarrow 0$ . Thus the generalized eigenvalue  $\lambda^* = 0$  in this case belongs to the residual spectrum of  $A$ . In either case,  $\lim_{j \rightarrow \infty} e_j$  is the corresponding generalized eigenvector that enlarges the trivial null space  $\mathcal{N}(\mathcal{L}_{\lambda^*})$  of the generalized eigenvalue  $\lambda^*$ . In fact in these two and the Dirac delta example of Sec. 5 of continuous and residual spectra, the generalized eigenfunctions arise as the limits of a sequence

of functions whose images under the respective  $\mathcal{L}_\lambda$  converge to 0; recall the definition of footnote 26. This observation generalizes to the dense extension  $\text{Multi}_|(X, Y)$  of  $\text{Map}(X, Y)$  as follows. If  $x \in \mathcal{D}_+$  is not a fixed point of  $f(\lambda; x) = x$ , but there is some  $n \in \mathbb{N}$  such that  $f^n(\lambda; x) = x$ , then the limit  $n \rightarrow \infty$  generates a multifunction at  $x$  as was the case with the delta function in the previous section and the various other examples that we have seen so far in the earlier sections.

One of the main goals of investigations on the spectrum of nonlinear operators is to find a set in the complex plane that has the usual desirable properties of the spectrum of a linear operator [Appell *et al.*, 2000]. In this case, the focus has been to find a suitable class of operators  $\mathcal{C}(X)$  with  $T \in \mathcal{C}(X)$ , such that the resolvent set is expressed as

$$\rho(T) = \{\lambda \in \mathbb{C} : (T_\lambda \text{ is } 1 : 1)(\text{Cl}(\mathcal{R}(T_\lambda)) = X) \\ \text{and } (T_\lambda^{-1} \in \mathcal{C}(X) \text{ on } \mathcal{R}(T_\lambda))\}$$

with the spectrum  $\sigma(T)$  being defined as the complement of this set. Among the classes  $\mathcal{C}(X)$  that have been considered, beside spaces of continuous functions  $C(X)$ , are linear boundedness  $B(X)$ , Frechet differentiability  $C^1(X)$ , Lipschitz continuity  $\text{Lip}(X)$ , and Granas quasiboundedness  $Q(x)$ , where  $\text{Lip}(X)$  specifically takes into account the nonlinearity of  $T$  to define

$$\|T\|_{\text{Lip}} = \sup_{x \neq y} \frac{\|T(x) - T(y)\|}{\|x - y\|}, \\ |T|_{\text{lip}} = \inf_{x \neq y} = \frac{\|T(x) - T(y)\|}{\|x - y\|} \tag{72}$$

that are plain generalizations of the corresponding norms of linear operators. Plots of  $f_\lambda^-(y) = \{x \in \mathcal{D}(f - \lambda) : (f - \lambda)x = y\}$  for the functions  $f : \mathbb{R} \rightarrow \mathbb{R}$

$$f_{\lambda a}(x) = \begin{cases} -1 - \lambda x, & x < -1 \\ (1 - \lambda)x, & -1 \leq x \leq 1 \\ 1 - \lambda x, & 1 < x \end{cases}$$

$$f_{\lambda b}(x) = \begin{cases} -\lambda x, & x < 1 \\ (1 - \lambda)x - 1, & 1 \leq x \leq 2 \\ 1 - \lambda x, & 2 < x \end{cases}$$

$$f_{\lambda c}(x) = \begin{cases} -\lambda x & x < 1 \\ \sqrt{x - 1} - \lambda x & 1 \leq x \end{cases}$$

$$f_{\lambda d}(x) = \begin{cases} (x - 1)^2 + 1 - \lambda x & 1 \leq x \leq 2 \\ (1 - \lambda)x & \text{otherwise} \end{cases}$$

$$f_{\lambda e}(x) = \tan^{-1}(x) - \lambda x$$



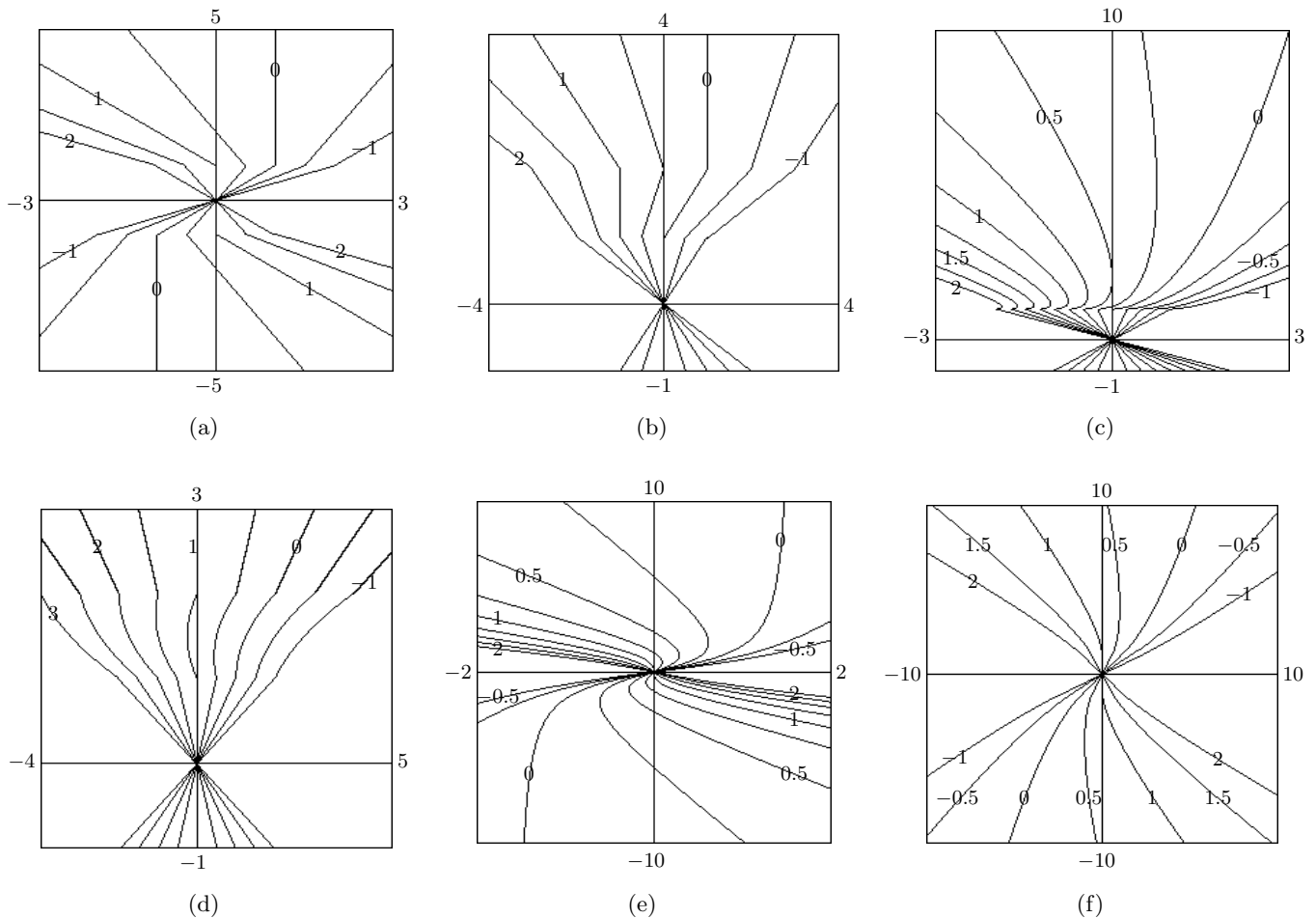


Fig. 19. Inverses of  $f_\lambda = f - \lambda$ . The  $\lambda$ -values are shown on the graphs.

$$f_\lambda(x) = \begin{cases} 1 - 2\sqrt{-x} - \lambda x, & x < -1 \\ (1 - \lambda)x, & -1 \leq x \leq 1 \\ 2\sqrt{x} - 1 - \lambda x, & 1 < x \end{cases}$$

taken from [Appell *et al.*, 2000] are shown in Fig. 19. It is easy to verify that the Lipschitz and linear upper and lower bounds of these maps are as in Table 4.

The point spectrum defined by

$$P\sigma(f) = \{\lambda \in \mathbb{C} : (f - \lambda)x = 0 \text{ for some } x \neq 0\}$$

is the simplest to calculate. Because of the special role played by the zero element 0 in generating the point spectrum in the linear case, the bounds  $m\|x\| \leq \|\mathcal{L}x\| \leq M\|x\|$  together with  $\mathcal{L}x = \lambda x$  imply  $\text{Cl}(P\sigma(\mathcal{L})) = [[\|\mathcal{L}\|_b, \|\mathcal{L}\|_B]]$  — where the subscripts denote the lower and upper bounds in Eq. (72) and which sometimes is taken to be a descriptor of the point spectrum of a nonlinear operator — as can be seen in Table 5 and verified from Fig. 19. The remainder of the spectrum, as

the complement of the resolvent set, is more difficult to find. Here the convenient characterization of the resolvent of a continuous linear operator as the set of all sufficiently large  $\lambda$  that satisfy  $|\lambda| > M$  is of little significance as, unlike for a linear operator, the non-existence of an inverse is not just due to the set  $\{f^{-1}(0)\}$  which happens to be the only way a linear map can fail to be injective. Thus the map defined piecewise as  $\alpha + 2(1 - \alpha)x$  for  $0 \leq x < 1/2$  and  $2(1 - x)$  for  $1/2 \leq x \leq 1$ , with  $0 < \alpha < 1$ , is not invertible on its range although  $\{f^{-1}(0)\} = 1$ . Comparing Fig. 19 and Table 4, it is seen that in cases (b)–(d), the intervals  $[[f|_b, \|f\|_B]]$  are subsets of the  $\lambda$ -values for which the respective maps are not injective; this is to be compared with (a), (e) and (f) where the two sets are the same. Thus the linear bounds are not good indicators of the uniqueness properties of solution of nonlinear equations for which the Lipschitzian bounds are seen to be more appropriate.

Table 4. Bounds on the functions of Fig. 19.

Function	$ f _b$	$\ f\ _B$	$ f _{lip}$	$\ f\ _{Lip}$
$f_a$	0	1	0	1
$f_b$	0	1/2	0	1
$f_c$	0	1/2	0	$\infty$
$f_d$	$2(\sqrt{2} - 1)$	$\infty$	0	2
$f_e$	0	1	0	1
$f_f$	0	1	0	1

Table 5. Lipschitzian and point spectra of the functions of Fig. 19.

Functions	$\sigma_{Lip}(f)$	$P\sigma(f)$
$f_a$	[0, 1]	(0, 1)
$f_b$	[0, 1]	[0, 1/2]
$f_c$	[0, $\infty$ )	[0, 1/2]
$f_d$	[0, 2]	$[2(\sqrt{2} - 1), 1]$
$f_e$	[0, 1]	(0, 1)
$f_f$	[0, 1]	(0, 1)

In view of the above, we may draw the following conclusions. If we choose to work in the space of multifunctions  $\text{Multi}(X, \mathcal{T})$ , with  $\mathcal{T}$  the topology of pointwise biconvergence, when all functional relations are (multi)invertible on their ranges, we may make the following definition for the net of functions  $f(\lambda; x)$  satisfying  $f(\lambda; x) = x$ .

**Definition 6.1.** Let  $f(\lambda; \cdot) \in \text{Multi}(X, \mathcal{T})$  be a function. The resolvent set of  $f$  is given by

$$\rho(f) = \{ \lambda : (f(\lambda; \cdot))^{-1} \in \text{Map}(X, \mathcal{T}) \wedge (\text{Cl}(\mathcal{R}(f(\lambda; \cdot))) = X) \},$$

and any  $\lambda$  not in  $\rho$  is in the spectrum of  $f$ .

Thus apart from multifunctions,  $\lambda \in \sigma(f)$  also generates functions on the boundary of functional and non-functional relations in  $\text{Multi}(X, \mathcal{T})$ . While it is possible to classify the spectrum into point, continuous and residual subsets, as in the linear case, it is more meaningful for nonlinear operators to consider  $\lambda$  as being either in the *boundary spectrum*  $\text{Bdy}(\sigma(f))$  or in the *interior spectrum*  $\text{Int}(\sigma(f))$ , depending on whether or not the multifunction  $f(\lambda; \cdot)^-$  arises as the graphical limit of a net of functions in either  $\rho(f)$  or  $R\sigma(f)$ . This is suggested by the spectra arising from the second row of Table 1 (injective  $\mathcal{L}_\lambda$  and discontinuous  $\mathcal{L}_\lambda^{-1}$ ) that lies sandwiched in the  $\lambda$ -plane between the two components arising from the first and third rows, see [Naylor & Sell, 1971, Sec. 6.6], for example. According to this simple scheme, the spectral set is a closed set with its boundary and interior belonging to  $\text{Bdy}(\sigma(f))$  and  $\text{Int}(\sigma(f))$ , respectively. Table 6 shows this division for the examples in Fig. 19. Because 0 is no more significant than any other point in the domain of a nonlinear map in inducing non-injectivity, the division of the spectrum into the traditional sets would be as shown in Table 6; compare also with the conventional linear point spectrum of Table 5. In this nonlinear classification, the point spectrum consists of any  $\lambda$  for which the inverse  $f(\lambda; \cdot)^-$  is set-valued, irrespective of whether this is produced at 0 or not, while the continuous and residual spectra together comprise the boundary spectrum. Thus a  $\lambda$  can be both at the point and the continuous or residual spectra which need not be disjoint. The continuous and residual spectra are included in the boundary spectrum which may also contain parts of the point spectrum.

**Example 6.1.** To see how these concepts apply to linear mappings, consider the equation  $(D - \lambda)y(x) = r(x)$  where  $D = d/dx$  is the differential

Table 6. Nonlinear spectra of functions of Fig. 19. Compare the present point spectra with the usual linear spectra of Table 5.

Function	$\text{Int}(\sigma(f))$	$\text{Bdy}(\sigma(f))$	$P\sigma(f)$	$C\sigma(f)$	$R\sigma(f)$
$f_a$	(0, 1)	{0, 1}	[0, 1]	{1}	{0}
$f_b$	(0, 1)	{0, 1}	[0, 1]	{1}	{0}
$f_c$	(0, $\infty$ )	{0}	[0, $\infty$ )	{0}	$\emptyset$
$f_d$	(0, 2)	{0, 2}	(0, 2)	{0, 2}	$\emptyset$
$f_e$	(0, 1)	{0, 1}	(0, 1)	{1}	{0}
$f_f$	(0, 1)	{0, 1}	(0, 1)	{0, 1}	$\emptyset$

operator on  $L^2[0, \infty)$ , and let  $\lambda$  be real. For  $\lambda \neq 0$ , the unique solution of this equation in  $L^2[0, \infty)$ , is

$$y(x) = \begin{cases} e^{\lambda x} \left( y(0) + \int_0^x e^{-\lambda x'} r(x') dx' \right), & \lambda < 0 \\ e^{\lambda x} \left( y(0) - \int_x^\infty e^{-\lambda x'} r(x') dx' \right) & \lambda > 0 \end{cases}$$

showing that for  $\lambda > 0$  the inverse is functional so that  $\lambda \in (0, \infty)$  belongs to the resolvent of  $D$ . However, when  $\lambda < 0$ , apart from the  $y = 0$  solution (since we are dealing with a linear problem, only  $r = 0$  is to be considered),  $e^{\lambda x}$  is also in  $L^2[0, \infty)$  so that all such  $\lambda$  are in the point spectrum of  $D$ . For  $\lambda = 0$  and  $r \neq 0$ , the two solutions are not necessarily equal unless  $\int_0^\infty r(x) = 0$ , so that the range  $\mathcal{R}(D - I)$  is a subspace of  $L^2[0, \infty)$ . To complete the problem, it is possible to show [Naylor & Sell, 1971] that  $0 \in C\sigma(D)$ , see Example 2.2; hence the continuous spectrum forms at the boundary of the functional solution for the resolvent- $\lambda$  and the multifunctional solution for the point spectrum. With a slight variation of problem to  $y(0) = 0$ , all  $\lambda < 0$  are in the resolvent set, while  $\lambda > 0$  the inverse is bounded but must satisfy  $y(0) = \int_0^\infty e^{-\lambda x} r(x) dx = 0$  so that  $\text{Cl}(\mathcal{R}(D - \lambda)) \neq L^2[0, \infty)$ . Hence  $\lambda > 0$  belong to the residual spectrum. The decomposition of the complex  $\lambda$ -plane for these and some other linear spectral problems taken from [Naylor & Sell, 1971] is shown in Fig. 20. In all cases, the spectrum due to the second row of Table 1 acts as a boundary between that arising from the first and third rows, which justifies our division of the spectrum for a nonlinear operator into the interior and boundary components. Compare with Example 2.2.

From the basic representation of the resolvent operator  $(\mathbf{1} - f)^{-1}$

$$\mathbf{1} + f + f^2 + \dots + f^i + \dots$$

in  $\text{Multi}(X)$ , if the iterates of  $f$  converge to a multifunction for some  $\lambda$ , then that  $\lambda$  must be in the spectrum of  $f$ , which means that the control parameter of a chaotic dynamical system is in its spectrum. Of course, the series can sum to a multi even otherwise: take  $f_\lambda(x)$  to be identically  $x$  with  $\lambda = 1$ , for example, to get  $1 \in P\sigma(f)$ . A comparison of Tables 1 and 5 reveal that in case (d), for example, 0 and 2 belong to the Lipschitz spectrum because although  $f_d^{-1}$  is not Lipschitz continuous,  $\|f\|_{\text{Lip}} = 2$ . It should also be noted that the boundary between the functional resolvent and multifunctional spectral set is formed

by the graphical convergence of a net of resolvent functions while the multifunctions in the interior of the spectral set evolve graphically independent of the functions in the resolvent. The chaotic states forming the boundary of the functional and multifunctional subsets of  $\text{Multi}(X)$  marks the transition from the less efficient functional state to the more efficient multifunctional one.

These arguments also suggest the following. The countably many outputs arising from the non-injectivity of  $f(\lambda; \cdot)$  corresponding to a given input can be interpreted to define *complexity because in a nonlinear system each of these possibilities constitute an experimental result in itself that may not be combined in any definite predetermined manner*. This is in sharp contrast to linear systems where a linear combination, governed by the initial conditions, always generate a unique end result; recall also the combination offered by the singular generalized eigenfunctions of neutron transport theory. This multiplicity of possibilities that have no definite combinatorial property is the basis of the diversity of nature, and is possibly responsible for Feigenbaum's "historical prejudice", [Feigenbaum, 1992], see Prelude 2. Thus *order* represented by the functional resolvent passes over to *complexity* of the countably multifunctional interior spectrum via the uncountably multifunctional boundary that is a prerequisite for *chaos*. We may now strengthen our hypothesis offered at the end of the previous section in terms of the examples of Figs. 19 and 20, that nature uses chaoticity as an intermediate step to the attainment of states that would otherwise be inaccessible to it. Well-posedness of a system is an extremely inefficient way of expressing a multitude of possibilities as this requires a different input for every possible output. Nature chooses to express its myriad manifestations through the multifunctional route leading either to averaging as in the delta function case or to a countable set of well-defined states, as in the examples of Fig. 19 corresponding to the interior spectrum. Of course it is no distraction that the multifunctional states arise respectively from  $f_\lambda$  and  $f_\lambda^-$  in these examples as  $f$  is a function on  $X$  that is under the influence of both  $f$  and its inverse. The functional resolvent is, for all practical purposes, only a tool in this structure of nature.

The equation  $f(x) = y$  is typically an input-output system in which the inverse images at a functional value  $y_0$  represents a set of input parameters leading to the same experimental output  $y_0$ ; this

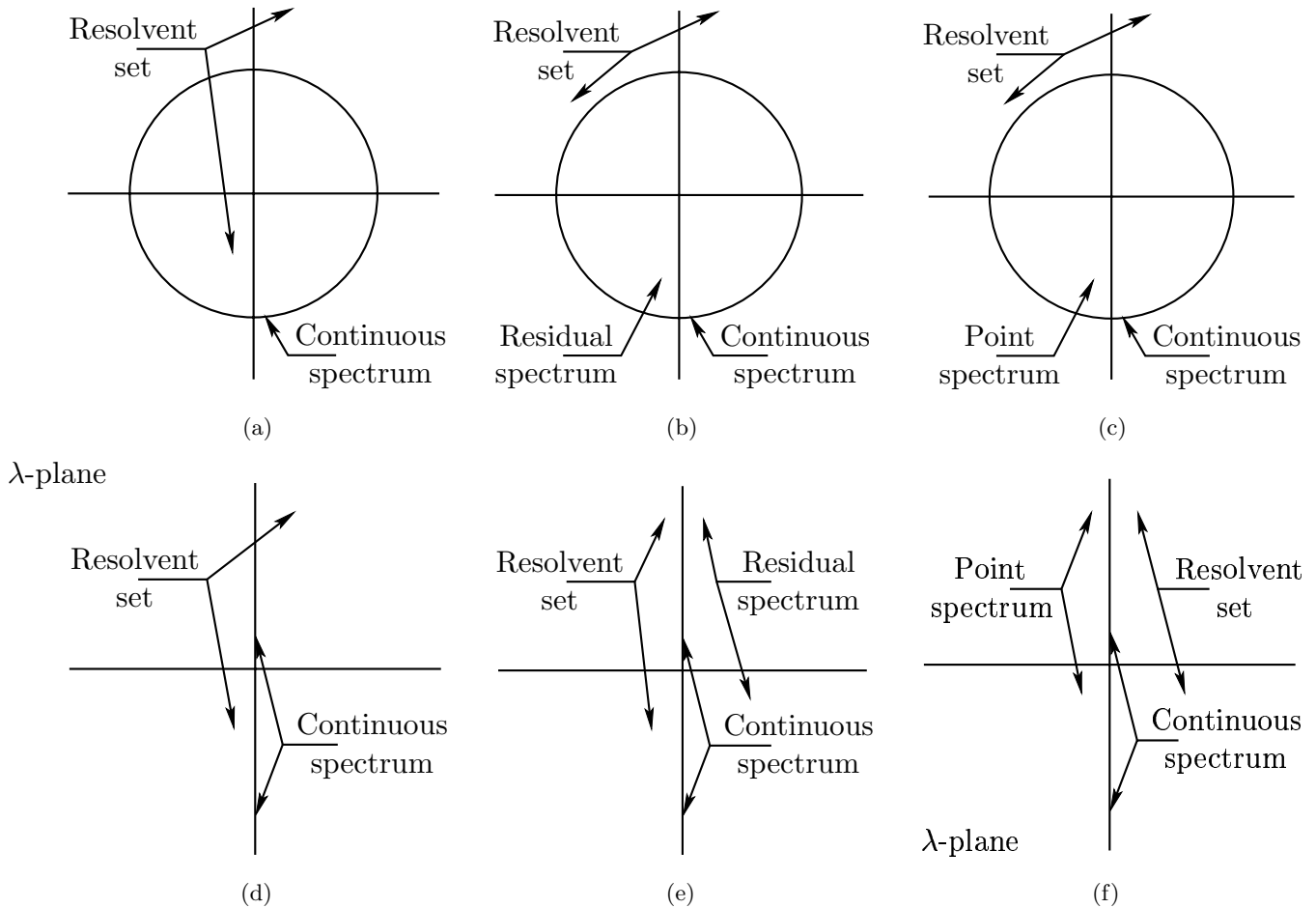


Fig. 20. Spectra of some linear operators in the complex  $\lambda$ -plane. (a) Left shift  $(\dots, x_{-1}, x_0, x_1, \dots) \mapsto (\dots, x_0, x_1, x_2, \dots)$  on  $l_2(-\infty, \infty)$ , (b) Right shift  $(x_0, x_1, x_2, \dots) \mapsto (0, x_0, x_1, \dots)$  on  $l_2[0, \infty)$ , (c) Left shift  $(x_0, x_1, x_2, \dots) \mapsto (x_1, x_2, x_3, \dots)$  on  $l_2[0, \infty)$  of sequence spaces, and (d)  $d/dx$  on  $L_2(-\infty, \infty)$  (e)  $d/dx$  on  $L_2[0, \infty)$  with  $y(0) = 0$  and (f)  $d/dx$  on  $L_2[0, \infty)$ . The residual spectrum in (b) and (e) arise from block (3–3) in Table 1, i.e.  $\mathcal{L}_\lambda$  is one-to-one and  $\mathcal{L}_\lambda^{-1}$  is bounded on non-dense domains in  $l_2[0, \infty)$  and  $L_2[0, \infty)$ , respectively. The continuous spectrum therefore marks the boundary between two functional states, as in (a) and (e), now with dense and non-dense domains of the inverse operator.

is stability characterized by a complete insensitivity of the output to changes in input. On the other hand, a continuous multifunction at  $x_0$  is a signal for a hypersensitivity to input because the output, which is a definite experimental quantity, is a choice from the possibly infinite set  $\{f(x_0)\}$  made by a choice function which represents the experiment at that particular point in time. Since there will always be finite differences in the experimental parameters when an experiment is repeated, the choice function (that is the experimental output) will select a point from  $\{f(x_0)\}$  that is representative of that experiment and which need not bear any definite relation to the previous values; this is instability and signals sensitivity to initial conditions. Such a state is of high entropy as the number of available states  $f_C(\{f(x_0)\})$  — where  $f_C$  is the choice function —

is larger than a functional state represented by the singleton  $\{f(x_0)\}$ .

### Epilogue

*The most passionate advocates of the new science go so far as to say that twentieth-century science will be remembered for just three things: relativity, quantum mechanics and chaos. Chaos, they contend, has become the century's third great revolution in the physical sciences. Like the first two revolutions, chaos cuts away at the tenets of Newton's physics. As one physicist put it: "Relativity eliminated the Newtonian illusion of absolute space and time; quantum theory eliminated the Newtonian dream of a controllable measurement process; and chaos eliminates the Laplacian fantasy of deterministic predictability." Of the three, the revolution in chaos applies to the*

universe we see and touch, to objects at human scale. . . . There has long been a feeling, not always expressed openly, that theoretical physics has strayed far from human intuition about the world. Whether this will prove to be fruitful heresy, or just plain heresy, no one knows. But some of those who thought that physics might be working its way into a corner now look to chaos as a way out.

[Gleick, 1987]

## Acknowledgments

It is a pleasure to thank the referees for recommending an enlarged Tutorial and Review revision of the original submission *Graphical Convergence, Chaos and Complexity*, and Professor Leon O. Chua for suggesting a pedagogically self-contained, jargonless version accessible to a wider audience for the present form of the paper. Financial assistance during the initial stages of this work from the National Board for Higher Mathematics is also acknowledged.

## References

- Alligood, K. T., Sauer, T. D. & Yorke, J. A. [1997] *Chaos, An Introduction to Dynamical Systems* (Springer-Verlag, NY).
- Appell, J., DePascale, E. & Vignoli, A. [2000] "A comparison of different spectra for nonlinear operators," *Nonlin. Anal.* **40**, 73–90.
- Brown, R. & Chua, L. O. [1996] "Clarifying chaos: Examples and counterexamples," *Int. J. Bifurcation and Chaos* **6**, 219–249.
- Campbell, S. I. & Mayer, C. D. [1979] *Generalized Inverses of Linear Transformations* (Pitman Publishing Ltd., London).
- Case, K. M. & Zweifel, P. F. [1967] *Linear Transport Theory* (Addison-Wesley, MA).
- de Souza, H. G. [1997] "Opening address," in *The Impact of Chaos on Science and Society*, eds. Grebogi, C. & Yorke, J. A. (United Nations University Press, Tokyo), pp. 384–386.
- Devaney, R. L. [1989] *An Introduction to Chaotic Dynamical Systems* (Addison-Wesley, CA).
- Falconer, K. [1990] *Fractal Geometry* (John Wiley, Chichester).
- Feigenbaum, M. [1992] "Foreword," *Chaos and Fractals: New Frontiers of Science* (Springer-Verlag, NY), pp. 1–7.
- Gallagher, R. & Appenzeller, T. [1999] "Beyond reductionism," *Science* **284**, p. 79.
- Gleick, J. [1987] *Chaos: The Amazing Science of the Unpredictable* (Viking, NY).
- Goldenfeld, N. & Kadanoff, L. P. [1999] "Simple lessons from complexity," *Science* **284**, 87–89.
- Korevaar, J. [1968] *Mathematical Methods*, Vol. 1 (Academic Press, NY).
- Naylor, A. W. & Sell, G. R. [1971] *Linear Operator Theory is Engineering and Science* Holt (Rinehart and Winston, NY).
- Peitgen, H.-O., Jurgens, H. & Saupe, D. [1992] *Chaos and Fractals: New Frontiers of Science* (Springer-Verlag, NY).
- Robinson, C. [1999] *Dynamical Systems: Stability, Symbolic Dynamics and Chaos* (CRC Press LLC, Boca Raton).
- Roman, P. [1975] *Some Modern Mathematics for Physicists and other Outsiders* (Pergamon Press, NY).
- Sengupta, A. [1995a] "A discretized spectral approximation in neutron transport theory. Some numerical considerations," *J. Stat. Phys.* **51**, 657–676.
- Sengupta, A. [1995b] "Full range solution of half-space neutron transport problem," *ZAMP* **46**, 40–60.
- Sengupta, A. [1997] "Multifunction and generalized inverse," *J. Inverse and Ill-Posed Problems* **5**, 265–285.
- Sengupta, A. & Ray, G. G. [2000] "A multifunctional extension of function spaces: Chaotic systems are maximally ill-posed," *J. Inverse and Ill-Posed Problems* **8**, 232–353.
- Stuart, A. M. & Humphries, A. R. [1996] *Dynamical Systems and Numerical Analysis* (Cambridge University Press).
- Tikhonov, A. N. & Arsenin, V. Y. [1977] *Solutions of Ill-Posed Problems* (V. H. Winston, Washington D.C.).
- Waldrop, M. M. [1992] *Complexity: The Emerging Science at the Edge of Order and Chaos* (Simon and Schuster).
- Willard, S. [1970] *General Topology* (Addison-Wesley, Reading, MA).
- Williams, M. M. R. [1971] *Mathematical Methods of Particle Transport Theory* (Butterworths, London).

## Appendix

This Appendix gives a brief overview of some aspects of topology that are necessary for a proper understanding of the concepts introduced in this work.

### A.1. Convergence in Topological Spaces: Sequence, Net and Filter

In the theory of convergence in topological spaces, *countability* plays an important role. To understand the significance of this concept, some preliminaries are needed.

The notion of a basis, or base, is a familiar one in analysis: a base is a subcollection of a set which may be used to construct, in a specified manner, any element of the set. This simplifies the statement of a problem since a smaller number of elements of the base can be used to generate the larger class of every element of the set. This philosophy finds application in topological spaces as follows.

Among the three properties (N1) – (N3) of the neighborhood system  $\mathcal{N}_x$  of Tutorial 4, (N1) and (N2) are basic in the sense that the resulting subcollection of  $\mathcal{N}_x$  can be used to generate the full system by applying (N3); this *basic neighborhood system*, or *neighborhood (local) base*  $\mathcal{B}_x$  at  $x$ , is characterized by

- (NB1)  $x$  belongs to each member  $B$  of  $\mathcal{B}_x$ .
- (NB2) The intersection of any two members of  $\mathcal{B}_x$  contains another member of  $\mathcal{B}_x$ :  $B_1, B_2 \in \mathcal{B}_x \Rightarrow (\exists B \in \mathcal{B}_x : B \subseteq B_1 \cap B_2)$ .

Formally, compare Eq. (18),

**Definition A.1.1.** A neighborhood (local) base  $\mathcal{B}_x$  at  $x$  in a topological space  $(X, \mathcal{U})$  is a subcollection of the neighborhood system  $\mathcal{N}_x$  having the property that each  $N \in \mathcal{N}_x$  contains some member of  $\mathcal{B}_x$ . Thus

$$\mathcal{B}_x \stackrel{\text{def}}{=} \{B \in \mathcal{N}_x : x \in B \subseteq N \text{ for each } N \in \mathcal{N}_x\} \tag{A.1}$$

determines the full neighborhood system

$$\mathcal{N}_x = \{N \subseteq X : x \in B \subseteq N \text{ for some } B \in \mathcal{B}_x\} \tag{A.2}$$

reciprocally as all supersets of the basic elements.

The entire neighborhood system  $\mathcal{N}_x$ , which is recovered from the base by forming all supersets of the basic neighborhoods, is trivially a local base at  $x$ ; non-trivial examples are given below.

The second example of a base, consisting as usual of a subcollection of a given collection, is the topological base  ${}_{\top}\mathcal{B}$  that allows the specification of the topology on a set  $X$  in terms of a smaller collection of open sets.

**Definition A.1.2.** A base  ${}_{\top}\mathcal{B}$  in a topological space  $(X, \mathcal{U})$  is a subcollection of the topology  $\mathcal{U}$  having the property that each  $U \in \mathcal{U}$  contains some member of  ${}_{\top}\mathcal{B}$ . Thus

$${}_{\top}\mathcal{B} \stackrel{\text{def}}{=} \{B \in \mathcal{U} : B \subseteq U \text{ for each } U \in \mathcal{U}\} \tag{A.3}$$

determines reciprocally the topology  $\mathcal{U}$  as

$$\mathcal{U} = \left\{ U \subseteq X : U = \bigcup_{B \in {}_{\top}\mathcal{B}} B \right\}. \tag{A.4}$$

This means that the topology on  $X$  can be reconstructed from the base by taking all possible unions of members of the base, and a collection of subsets of a set  $X$  is a topological base iff Eq. (A.4) of arbitrary unions of elements of  ${}_{\top}\mathcal{B}$  generates a topology on  $X$ . This topology, which is the coarsest (that is the smallest) that contains  ${}_{\top}\mathcal{B}$ , is obviously closed under finite intersections. Since the open set  $\text{Int}(N)$  is a neighborhood of  $x$  whenever  $N$  is, Eq. (A.2) and the definition Eq. (17) of  $\mathcal{N}_x$  implies that *the open neighborhood system of any point in a topological space is an example of a neighborhood base at that point*, an observation that has often led, together with Eq. (A.3), to the use of the term “neighborhood” as a synonym for “non-empty open set”. The distinction between the two however is significant as neighborhoods need not necessarily be open sets; thus while not necessary, it is clearly sufficient for the local basic sets  $B$  to be open in Eqs. (A.1) and (A.2). If Eq. (A.2) holds for every  $x \in N$ , then the resulting  $\mathcal{N}_x$  reduces to the topology induced by the open basic neighborhood system  $\mathcal{B}_x$  as given by Eq. (18).

In order to check if a collection of subsets  ${}_{\top}\mathcal{B}$  of  $X$  qualifies to be a basis, it is not necessary to verify properties (T1)–(T3) of Tutorial 4 for the class (A.4) generated by it because of the properties (TB1) and (TB2) below whose strong affinity to (NB1) and (NB2) is formalized in Theorem A.1.1.

**Theorem A.1.1.** A collection  ${}_{\top}\mathcal{B}$  of subsets of  $X$  is a topological basis on  $X$  iff

(TB1)  $X = \bigcup_{B \in {}_{\top}\mathcal{B}} B$ . Thus each  $x \in X$  must belong to some  $B \in {}_{\top}\mathcal{B}$  which implies the existence of a local base at each point  $x \in X$ .

(TB2) The intersection of any two members  $B_1$  and  $B_2$  of  ${}_{\top}\mathcal{B}$  with  $x \in B_1 \cap B_2$  contains another member of  ${}_{\top}\mathcal{B}$ :  $(B_1, B_2 \in {}_{\top}\mathcal{B}) \wedge (x \in B_1 \cap B_2) \Rightarrow (\exists B \in {}_{\top}\mathcal{B} : x \in B \subseteq B_1 \cap B_2)$ .

This theorem, together with Eq. (A.4) ensures that a given collection of subsets of a set  $X$  satisfying (TB1) and (TB2) induces *some* topology on  $X$ ; compared to this is the result that *any* collection of subsets of a set  $X$  is a *subbasis* for some topology on  $X$ . If  $X$ , however, already has a topology  $\mathcal{U}$  imposed on it, then Eq. (A.3)

must also be satisfied in order that the topology generated by  $\tau\mathcal{B}$  is indeed  $\mathcal{U}$ . The next theorem connects the two types of bases of Definitions A.1.1 and A.1.2 by asserting that although a local base of a space need not consist of open sets and a topological base need not have any reference to a point of  $X$ , any subcollection of the base containing a point is a local base at that point.

**Theorem A.1.2.** *A collection of open sets  $\tau\mathcal{B}$  is a base for a topological space  $(X, \mathcal{U})$  iff for each  $x \in X$ , the subcollection*

$$\mathcal{B}_x = \{B \in \mathcal{U} : x \in B \in \tau\mathcal{B}\} \tag{A.5}$$

*of basic sets containing  $x$  is a local base at  $x$ .*

*Proof. Necessity.* Let  $\tau\mathcal{B}$  be a base of  $(X, \mathcal{U})$  and  $N$  be a neighborhood of  $x$ , so that  $x \in U \subseteq N$  for some open set  $U = \bigcup_{B \in \tau\mathcal{B}} B$  and basic open sets  $B$ . Hence  $x \in B \subseteq N$  shows, from Eq. (A.1), that  $B \in \mathcal{B}_x$  is a local basic set at  $x$ .

*Sufficiency.* If  $U$  is an open set of  $X$  containing  $x$ , then the definition of local base Eq. (A.1) requires  $x \in B_x \subseteq U$  for some subcollection of basic sets  $B_x$  in  $\mathcal{B}_x$ ; hence  $U = \bigcup_{x \in U} B_x$ . By Eq. (A.4) therefore,  $\tau\mathcal{B}$  is a topological base for  $X$ . ■

Because the basic sets are open, (TB2) of Theorem A.1.1 leads to the following physically appealing paraphrase of Theorem A.1.2.

**Corollary.** *A collection  $\tau\mathcal{B}$  of open sets of  $(X, \mathcal{U})$  is a topological base that generates  $\mathcal{U}$  iff for each open set  $U$  of  $X$  and each  $x \in U$  there is an open set  $B \in \tau\mathcal{B}$  such that  $x \in B \subseteq U$ ; that is iff*

$$x \in U \in \mathcal{U} \Rightarrow (\exists B \in \tau\mathcal{B} : x \in B \subseteq U).$$

**Example A.1.1.** Some examples of local bases in  $\mathbb{R}$  are intervals of the type  $(x-\varepsilon, x+\varepsilon)$ ,  $[x-\varepsilon, x+\varepsilon]$  for real  $\varepsilon$ ,  $(x-q, x+q)$  for rational  $q$ ,  $(x-1/n, x+1/n)$  for  $n \in \mathbb{Z}_+$ , while for a metrizable space with the topology induced by a metric  $d$ , each of the following is a local base at  $x \in X$ :  $B_\varepsilon(x; d) := \{y \in X : d(x, y) < \varepsilon\}$  and  $D_\varepsilon(x; d) := \{y \in X : d(x, y) \leq \varepsilon\}$  for  $\varepsilon > 0$ ,  $B_q(x; d)$  for  $\mathbb{Q} \ni q > 0$  and  $B_{1/n}(x; d)$  for  $n \in \mathbb{Z}_+$ . In  $\mathbb{R}^2$ , two neighborhood bases at any  $x \in \mathbb{R}^2$  are the disks centered at  $x$  and the set of all squares at  $x$  with sides parallel to the axes. Although these bases have no elements in common, they are nevertheless equivalent in the sense that they both generate the same (usual) topology in

$\mathbb{R}^2$ . Of course, the entire neighborhood system at any point of a topological space is itself a (less useful) local base at that point. By Theorem A.1.2,  $B_\varepsilon(x; d)$ ,  $D_\varepsilon(x; d)$ ,  $\varepsilon > 0$ ,  $B_q(x; d)$ ,  $\mathbb{Q} \ni q > 0$  and  $B_{1/n}(x; d)$ ,  $n \in \mathbb{Z}_+$ , for all  $x \in X$  are examples of bases in a metrizable space with topology induced by a metric  $d$ .

In terms of local bases and bases, it is now possible to formulate the notions of first and second countability as follows.

**Definition A.1.3.** A topological space is first countable if each  $x \in X$  has some countable neighborhood base, and is second countable if it has a countable base.

Every metrizable space  $(X, d)$  is first countable as both  $\{B(x, q)\}_{\mathbb{Q} \ni q > 0}$  and  $\{B(x, 1/n)\}_{n \in \mathbb{Z}_+}$  are examples of countable neighborhood bases at any  $x \in (X, d)$ ; hence  $\mathbb{R}^n$  is first countable. It should be clear that although every second countable space is first countable, *only a countable first countable space can be second countable*, and a common example of an uncountable first countable space that is also second countable is provided by  $\mathbb{R}^n$ . Metrizable spaces need not be second countable: any uncountable set having the discrete topology is as an example.

**Example A.1.2.** The following is an important example of a space that is not first countable as it is needed for our pointwise biconvergence of Sec. 3. Let  $\text{Map}(X, Y)$  be the set of all functions between the uncountable spaces  $(X, \mathcal{U})$  and  $(Y, \mathcal{V})$ . Given any integer  $I \geq 1$ , and any finite collection of points  $(x_i)_{i=1}^I$  of  $X$  and of open sets  $(V_i)_{i=1}^I$  in  $Y$ , let

$$B((x_i)_{i=1}^I; (V_i)_{i=1}^I) = \{g \in \text{Map}(X, Y) : (g(x_i) \in V_i)(i = 1, 2, \dots, I)\} \tag{A.6}$$

be the functions in  $\text{Map}(X, Y)$  whose graphs pass through each of the sets  $(V_i)_{i=1}^I$  at  $(x_i)_{i=1}^I$ , and let  $\tau\mathcal{B}$  be the collection of all such subsets of  $\text{Map}(X, Y)$  for every choice of  $I$ ,  $(x_i)_{i=1}^I$ , and  $(V_i)_{i=1}^I$ . The existence of a unique topology  $\mathcal{T}$  — the *topology of pointwise convergence* on  $\text{Map}(X, Y)$  — that is generated by the open sets  $B$  of the collection  $\tau\mathcal{B}$  now follows because

(TB1) is satisfied: For any  $f \in \text{Map}(X, Y)$  there must be some  $x \in X$  and a corresponding  $V \subseteq Y$  such that  $f(x) \in V$ , and

(TB2) is satisfied because

$$\begin{aligned} & B((s_i)_{i=1}^I; (V_i)_{i=1}^I) \cap B((t_j)_{j=1}^J; (W_j)_{j=1}^J) \\ &= B((s_i)_{i=1}^I, (t_j)_{j=1}^J; (V_i)_{i=1}^I, (W_j)_{j=1}^J) \end{aligned}$$

implies that a function simultaneously belonging to the two open sets on the left must pass through each of the points defining the open set on the right.

We now demonstrate that  $(\text{Map}(X, Y), \mathcal{T})$  is not first countable by verifying that it is not possible to have a countable local base at any  $f \in \text{Map}(X, Y)$ . If this is not indeed true, let  $B_f^I((x_i)_{i=1}^I; (V_i)_{i=1}^I) = \{g \in \text{Map}(X, Y) : (g(x_i) \in V_i)_{i=1}^I\}$ , which denotes those members of  ${}_{\mathcal{T}}\mathcal{B}$  that contain  $f$  with  $V_i$  an open neighborhood of  $f(x_i)$  in  $Y$ , be a countable local base at  $f$ , see Theorem A.1.2. Since  $X$  is uncountable, it is now possible to choose some  $x^* \in X$  different from any of the  $(x_i)_{i=1}^I$  (for example, let  $x^* \in \mathbb{R}$  be an irrational for rational  $(x_i)_{i=1}^I$ ), and let  $f(x^*) \in V^*$  where  $V^*$  is an open neighborhood of  $f(x^*)$ . Then  $B(x^*; V^*)$  is an open set in  $\text{Map}(X, Y)$  containing  $f$ ; hence from the definition of the local base, Eq. (A.1), or equivalently from the Corollary to Theorem A.1.2, there exists some (countable)  $I \in \mathbb{N}$  such that  $f \in B^I \subseteq B(x^*; V^*)$ . However,

$$f^*(x) = \begin{cases} y_i \in V_i, & \text{if } x = x_i, \text{ and } 1 \leq i \leq I \\ y^* \notin V^*, & \text{if } x = x^* \\ \text{arbitrary,} & \text{otherwise} \end{cases}$$

is a simple example of a function on  $X$  that is in  $B^I$  (as it is immaterial as to what values the function takes at points other than those defining  $B^I$ ), but not in  $B(x^*; V^*)$ . From this it follows that a *sufficient condition for the topology of pointwise convergence to be first countable is that  $X$  be countable.*

Even though it is not first countable,  $(\text{Map}(X, Y), \mathcal{T})$  is a Hausdorff space when  $Y$  is Hausdorff. Indeed, if  $f, g \in (\text{Map}(X, Y), \mathcal{T})$  with  $f \neq g$ , then  $f(x) \neq g(x)$  for some  $x \in X$ . But then as  $Y$  is Hausdorff, it is possible to choose disjoint open intervals  $V_f$  and  $V_g$  at  $f(x)$  and  $g(x)$  respectively.

With this background on first and second countability, it is now possible to go back to the question of nets, filters and sequences. Technically, a sequence on a set  $X$  is a map  $x : \mathbb{N} \rightarrow X$  from the set of natural numbers to  $X$ ; instead of denoting this in the usual functional manner of  $x(i)$  with  $i \in \mathbb{N}$ , it is the standard practice to use the notation  $(x_i)_{i \in \mathbb{N}}$  for the terms of a sequence. However, if

the space  $(X, \mathcal{U})$  is not first countable (and as seen above this is not a rare situation), it is not difficult to realize that sequences are inadequate to describe convergence in  $X$  simply because it can have only countably many values whereas the space may require uncountably many neighborhoods to completely define the neighborhood system at a point. The resulting uncountable generalizations of a sequence in the form of *nets* and *filters* is achieved through a corresponding generalization of the index set  $\mathbb{N}$  to the directed set  $\mathbb{D}$ .

**Definition A.1.4.** A directed set  $\mathbb{D}$  is a preordered set for which the order  $\preceq$ , known as a direction of  $\mathbb{D}$ , satisfies

- (a)  $\alpha \in \mathbb{D} \Rightarrow \alpha \preceq \alpha$  (that is  $\preceq$  is reflexive).
- (b)  $\alpha, \beta, \gamma \in \mathbb{D}$  such that  $(\alpha \preceq \beta \wedge \beta \preceq \gamma) \Rightarrow \alpha \preceq \gamma$  (that is  $\preceq$  is transitive).
- (c)  $\alpha, \beta \in \mathbb{D} \Rightarrow \exists \gamma \in \mathbb{D}$  such that  $(\alpha \preceq \gamma \wedge \beta \preceq \gamma)$ .

While the first two properties are obvious enough, the third which replaces antisymmetry, ensures that for any finite number of elements of the directed set, there is always a successor (upper bound). Examples of directed sets can be both straight forward, as any totally ordered set like  $\mathbb{N}, \mathbb{R}, \mathbb{Q}$ , or  $\mathbb{Z}$  and all subsets of a set  $X$  under the superset or subset relation (that is  $(\mathcal{P}(X), \supseteq)$  or  $(\mathcal{P}(X), \subseteq)$  that are directed by their usual ordering, and not quite so obvious as the following examples which are significantly useful in dealing with convergence questions in topological spaces, amply illustrate.

The neighborhood system

$${}_{\mathbb{D}}N = \{N : N \in \mathcal{N}_x\}$$

at a point  $x \in X$ , directed by the reverse inclusion direction  $\preceq$  defined as

$$M \preceq N \Leftrightarrow N \subseteq M \quad \text{for } M, N \in \mathcal{N}_x, \quad (\text{A.7})$$

is a fundamental example of a *natural direction of  $\mathcal{N}_x$* . In fact while reflexivity and transitivity are clearly obvious, (c) follows because for any  $M, N \in \mathcal{N}_x$ ,  $M \preceq M \cap N$  and  $N \preceq M \cap N$ . Of course, this direction is not a total ordering on  $\mathcal{N}_x$ . A more naturally useful directed set in convergence theory is

$${}_{\mathbb{D}}N_t = \{(N, t) : (N \in \mathcal{N}_x)(t \in N)\} \quad (\text{A.8})$$

under its *natural direction*

$$(M, s) \preceq (N, t) \Leftrightarrow N \subseteq M \quad \text{for } M, N \in \mathcal{N}_x; \quad (\text{A.9})$$

${}_{\mathbb{D}}N_t$  is more useful than  ${}_{\mathbb{D}}N$  because, unlike the latter,  ${}_{\mathbb{D}}N_t$  does not require a simultaneous choice



of points from every  $N \in \mathcal{N}_x$  that implicitly involves a simultaneous application of the Axiom of Choice; see Example A.1.3 below. The general indexed variation

$$\mathbb{D}N_\beta = \{(N, \beta) : (N \in \mathcal{N}_x)(\beta \in \mathbb{D})(x_\beta \in N)\} \tag{A.10}$$

of Eq. (A.8), with natural direction

$$(M, \alpha) \leq (N, \beta) \Leftrightarrow (\alpha \preceq \beta) \wedge (N \subseteq M), \tag{A.11}$$

often proves useful in applications as will be clear from the proofs of Theorems A.1.3 and A.1.4.

**Definition A.1.5 (Net).** Let  $X$  be any set and  $\mathbb{D}$  a directed set. A net  $\chi : \mathbb{D} \rightarrow X$  in  $X$  is a function on the directed set  $\mathbb{D}$  with values in  $X$ .

A net, to be denoted as  $\chi(\alpha)$ ,  $\alpha \in \mathbb{D}$ , is therefore a function indexed by a directed set. We adopt the convention of denoting nets in the manner of functions and do not use the sequential notation  $\chi_\alpha$  that can also be found in the literature. Thus, while every sequence is a special type of net,  $\chi : \mathbb{Z} \rightarrow X$  is an example of a net that is not a sequence.

Convergence of sequences and nets are described most conveniently in terms of the notions of being *eventually in* and *frequently in* every neighborhood of points. We describe these concepts in terms of nets which apply to sequences with obvious modifications.

**Definition A.1.6.** A net  $\chi : \mathbb{D} \rightarrow X$  is said to be

- (a) Eventually in a subset  $A$  of  $X$  if its tail is eventually in  $A$ :  $(\exists \beta \in \mathbb{D}) : (\forall \gamma \succeq \beta)(\chi(\gamma) \in A)$ .
- (b) Frequently in a subset  $A$  of  $X$  if for any index  $\beta \in \mathbb{D}$ , there is a successor index  $\gamma \in \mathbb{D}$  such that  $\chi(\gamma)$  is in  $A$ :  $(\forall \beta \in \mathbb{D})(\exists \gamma \succeq \beta) : (\chi(\gamma) \in A)$ .

It is not difficult to appreciate that

- (i) A net eventually in a subset is also frequently in it but not conversely,
- (ii) A net eventually (respectively, frequently) in a subset cannot be frequently (respectively, eventually) in its complement.

With these notions of eventually in and frequently in, convergence characteristics of a net may be expressed as follows.

**Definition A.1.7.** A net  $\chi : \mathbb{D} \rightarrow X$  converges to  $x \in X$  if it is eventually in every neighborhood of  $x$ , that is

$$(\forall N \in \mathcal{N}_x)(\exists \mu \in \mathbb{D})(\chi(\nu \succeq \mu) \in N).$$

The point  $x$  is known as the limit of  $\chi$  and the collection of all limits of a net is the limit set

$$\lim(\chi) = \{x \in X : (\forall N \in \mathcal{N}_x)(\exists \mathbb{R}_\beta \in \text{Res}(\mathbb{D}))(\chi(\mathbb{R}_\beta) \subseteq N)\} \tag{A.12}$$

of  $\chi$ , with the set of residuals  $\text{Res}(\mathbb{D})$  in  $\mathbb{D}$  given by

$$\text{Res}(\mathbb{D}) = \{\mathbb{R}_\alpha \in \mathcal{P}(\mathbb{D}) : \mathbb{R}_\alpha = \{\beta \in \mathbb{D} \text{ for all } \beta \succeq \alpha \in \mathbb{D}\}\}. \tag{A.13}$$

The net adheres at  $x \in X$ <sup>27</sup> if it is frequently in every neighborhood of  $x$ , that is

$$((\forall N \in \mathcal{N}_x)(\forall \mu \in \mathbb{D}))((\exists \nu \succeq \mu) : \chi(\nu) \in N).$$

The point  $x$  is known as the adherent of  $\chi$  and the collection of all adherents of  $\chi$  is the adherent set of the net, which may be expressed in terms of the cofinal subset of  $\mathbb{D}$

$$\text{Cof}(\mathbb{D}) = \{\mathbb{C}_\alpha \in \mathcal{P}(\mathbb{D}) : \mathbb{C}_\alpha = \{\beta \in \mathbb{D} \text{ for some } \beta \succeq \alpha \in \mathbb{D}\}\} \tag{A.14}$$

(thus  $\mathbb{D}_\alpha$  is cofinal in  $\mathbb{D}$  iff it intersects every residual in  $\mathbb{D}$ ), as

$$\text{adh}(\chi) = \{x \in X : (\forall N \in \mathcal{N}_x)(\exists \mathbb{C}_\beta \in \text{Cof}(\mathbb{D}))(\chi(\mathbb{C}_\beta) \subseteq N)\}. \tag{A.15}$$

This recognizes, in keeping with the limit set, each subnet of a net to be a net in its own right, and is equivalent to

$$\text{adh}(\chi) = \{x \in X : (\forall N \in \mathcal{N}_x)(\forall \mathbb{R}_\alpha \in \text{Res}(\mathbb{D}))(\chi(\mathbb{R}_\alpha) \cap N \neq \emptyset)\}. \tag{A.16}$$

Intuitively, a sequence is eventually in a set  $A$  if it is always in it after a finite number of terms (of course, the concept of a *finite number of terms* is unavailable for nets; in this case the situation may be described by saying that a net is eventually in  $A$  if its *tail is in A*) and it is frequently in  $A$  if it always returns to  $A$  to leave it again. It can be shown that a net is eventually (resp. frequently) in a set iff it is not frequently (resp. eventually) in its complement.

The following examples illustrate graphically the role of a proper choice of the index set  $\mathbb{D}$  in the description of convergence.

<sup>27</sup>This is also known as a *cluster point*; we shall, however, use this new term exclusively in the sense of the elements of a derived set, see Definition 2.3.

**Example A.1.3.** (1) Let  $\gamma \in \mathbb{D}$ . The eventually constant net  $\chi(\delta) = x$  for  $\delta \succeq \gamma$  converges to  $x$ .  
 (2) Let  $\mathcal{N}_x$  be a neighborhood system at a point  $x$  in  $X$  and suppose that the net  $(\chi(N))_{N \in \mathcal{N}_x}$  is defined by

$$\chi(M) \stackrel{\text{def}}{=} s \in M; \tag{A.17}$$

here the directed index set  ${}_{\mathbb{D}}N$  is ordered by the natural direction (A.7) of  $\mathcal{N}_x$ . Then  $\chi(N) \rightarrow x$  because given any  $x$ -neighborhood  $M \in {}_{\mathbb{D}}N$ , it follows from

$$M \preceq N \in {}_{\mathbb{D}}N \Rightarrow \chi(N) = t \in N \subseteq M \tag{A.18}$$

that a point in any subset of  $M$  is also in  $M$ ;  $\chi(N)$  is therefore eventually in every neighborhood of  $x$ .

(3) This slightly more general form of the previous example provides a link between the complimentary concepts of nets and filters that is considered below. For a point  $x \in X$ , and  $M, N \in \mathcal{N}_x$  with the corresponding directed set  ${}_{\mathbb{D}}M_s$  of Eq. (A.8) ordered by its natural order (A.9), the net

$$\chi(M, s) \stackrel{\text{def}}{=} s \tag{A.19}$$

converges to  $x$  because, as in the previous example, for any given  $(M, s) \in {}_{\mathbb{D}}N_s$ , it follows from

$$(M, s) \preceq (N, t) \in {}_{\mathbb{D}}M_s \Rightarrow \chi(N, t) = t \in N \subseteq M \tag{A.20}$$

that  $\chi(N, t)$  is eventually in every neighborhood  $M$  of  $x$ . The significance of the directed set  ${}_{\mathbb{D}}N_t$  of Eq. (A.8), as compared to  ${}_{\mathbb{D}}N$ , is evident from the net that it induces *without using the Axiom of Choice*: For a subset  $A$  of  $X$ , the net  $\chi(N, t) = t \in A$  indexed by the directed set

$${}_{\mathbb{D}}N_t = \{(N, t) : (N \in \mathcal{N}_x)(t \in N \cap A)\} \tag{A.21}$$

under the direction of Eq. (A.9), converges to  $x \in X$  with all such  $x$  defining the closure  $\text{Cl}(A)$  of  $A$ . Furthermore taking the directed set to be

$${}_{\mathbb{D}}N_t = \{(N, t) : (N \in \mathcal{N}_x)(t \in N \cap A - \{x\})\} \tag{A.22}$$

which, unlike Eq. (A.21), excludes the point  $x$  that may or may not be in the subset  $A$  of  $X$ , induces the net  $\chi(N, t) = t \in A - \{x\}$  converging to  $x \in X$ , with the set of all such  $x$  yielding the derived set  $\text{Der}(A)$  of  $A$ . In contrast, Eq. (A.21) also includes the isolated points  $t = x$  of  $A$  so as to generate its closure. Observe how neighborhoods of a point, which define convergence of nets and filters in a topological space  $X$ , double up here as index sets

to yield a self-consistent tool for the description of convergence.

As compared with sequences where, the index set is restricted to positive integers, the considerable freedom in the choice of directed sets as is abundantly borne out by the two preceding examples, is not without its associated drawbacks. Thus as a trade-off, the wide range of choice of the directed sets may imply that induction methods, so common in the analysis of sequences, need no longer apply to arbitrary nets.

(4) The non-convergent nets (actually these are sequences)

- (a)  $(1, -1, 1, -1, \dots)$  adheres at 1 and  $-1$  and
- (b)  $x_n = \begin{cases} n, & \text{if } n \text{ is odd} \\ 1 - 1/(1 + n), & \text{if } n \text{ is even} \end{cases}$

adheres at 1 for its even terms, but is unbounded in the odd terms.

A converging sequence or net is also adhering but, as examples (4) show, the converse is false. Nevertheless it is true, as again is evident from examples (4), that in a first countable space where sequences suffice, a sequence  $(x_n)$  adheres to  $x$  iff some subsequence  $(x_{n_m})_{m \in \mathbb{N}}$  of  $(x_n)$  converges to  $x$ . If the space is not first countable this has a corresponding equivalent formulation for nets with subnets replacing subsequences as follows.

Let  $(\chi(\alpha))_{\alpha \in \mathbb{D}}$  be a net. A *subnet* of  $\chi(\alpha)$  is the net  $\zeta(\beta) = \chi(\sigma(\beta))$ ,  $\beta \in \mathbb{E}$ , where  $\sigma : (\mathbb{E}, \preceq) \rightarrow (\mathbb{D}, \preceq)$  is a function that captures the essence of the subsequential mapping  $n \mapsto n_m$  in  $\mathbb{N}$  by satisfying

(SN1)  $\sigma$  is an increasing order-preserving function: it respects the order of  $\mathbb{E}$ :  $\sigma(\beta) \preceq \sigma(\beta')$  for every  $\beta \leq \beta' \in \mathbb{E}$ , and

(SN2) For every  $\alpha \in \mathbb{D}$  there exists a  $\beta \in \mathbb{E}$  such that  $\alpha \preceq \sigma(\beta)$ .

These generalize the essential properties of a subsequence in the sense that (1) Even though the index sets  $\mathbb{D}$  and  $\mathbb{E}$  may be different, it is necessary that the values of  $\mathbb{E}$  be contained in  $\mathbb{D}$ , and (2) There are arbitrarily large  $\alpha \in \mathbb{D}$  such that  $\chi(\alpha = \sigma(\beta))$  is a value of the subnet  $\zeta(\beta)$  for some  $\beta \in \mathbb{E}$ . Recalling the first of the order relations Eq. (38) on  $\text{Map}(X, Y)$ , we will denote a subnet  $\zeta$  of  $\chi$  by  $\zeta \preceq \chi$ .

We now consider the concept of filter on a set  $X$  that is very useful in visualizing the behavior of sequences and nets, and in fact filters constitute an alternate way of looking at convergence questions in topological spaces. A filter  $\mathcal{F}$  on a set  $X$

is a collection of *nonempty* subsets of  $X$  satisfying properties (F1) – (F3) below that are simply those of a neighborhood system  $\mathcal{N}_x$  without specification of the reference point  $x$ .

- (F1) The empty set  $\emptyset$  does not belong to  $\mathcal{F}$ ,
- (F2) The intersection of any two members of a filter is another member of the filter:  $F_1, F_2 \in \mathcal{F} \Rightarrow F_1 \cap F_2 \in \mathcal{F}$ ,
- (F3) Every superset of a member of a filter belongs to the filter:  $(F \in \mathcal{F}) \wedge (F \subseteq G) \Rightarrow G \in \mathcal{F}$ ; in particular  $X \in \mathcal{F}$ .

**Example A.1.4**

- (1) The *indiscrete filter* is the smallest filter on  $X$ .
- (2) The neighborhood system  $\mathcal{N}_x$  is the important *neighborhood filter at  $x$  on  $X$* , and any local base at  $x$  is also a filter-base for  $\mathcal{N}_x$ . In general for any subset  $A$  of  $X$ ,  $\{N \subseteq X : A \subseteq \text{Int}(N)\}$  is a filter on  $X$  at  $A$ .
- (3) All subsets of  $X$  containing a point  $x \in X$  is the *principal filter  ${}_F\mathcal{P}(x)$  on  $X$  at  $x$* . More generally, if  $\mathcal{F}$  consists of all supersets of a *nonempty* subset  $A$  of  $X$ , then  $\mathcal{F}$  is the *principal filter  ${}_F\mathcal{P}(A) = \{N \subseteq X : A \subseteq \text{Int}(N)\}$  at  $A$* . By adjoining the empty set to this filter give the  *$p$ -inclusion and  $A$ -inclusion topologies on  $X$ , respectively*. The single element sets  $\{\{x\}\}$  and  $\{A\}$  are particularly simple examples of filter-bases that generate the principal filters at  $x$  and  $A$ .
- (4) For an uncountable (resp. infinite) set  $X$ , all cocountable (resp. cofinite) subsets of  $X$  constitute the *cocountable (resp. cofinite or Frechet) filter on  $X$* . Again, adding to these filters the empty set gives the respective topologies.

Like the topological and local bases  ${}_T\mathcal{B}$  and  $\mathcal{B}_x$  respectively, a subclass of  $\mathcal{F}$  may be used to define a filter-base  ${}_F\mathcal{B}$  that in turn generate  $\mathcal{F}$  on  $X$ , just as it is possible to define the concepts of limit and adherence sets for a filter to parallel those for nets that follow straightforwardly from Definition A.1.7, taken with Definition A.1.11.

**Definition A.1.8.** Let  $(X, \mathcal{T})$  be a topological

---

<sup>28</sup>The restatement

$$\mathcal{F} \rightarrow x \Leftrightarrow \mathcal{N}_x \subseteq \mathcal{F} \tag{A.25}$$

of Eq. (A.23) that follows from (F3), and sometimes taken as the definition of convergence of a filter, is significant as it ties up the algebraic filter with the topological neighborhood system to produce the filter theory of convergence in topological spaces. From the defining properties of  $\mathcal{F}$  it follows that for each  $x \in X$ ,  $\mathcal{N}_x$  is the coarsest (that is smallest) filter on  $X$  that converges to  $x$ .

space and  $\mathcal{F}$  a filter on  $X$ . Then

$$\lim(\mathcal{F}) = \{x \in X : (\forall N \in \mathcal{N}_x)(\exists F \in \mathcal{F})(F \subseteq N)\} \tag{A.23}$$

and

$$\text{adh}(\mathcal{F}) = \{x \in X : (\forall N \in \mathcal{N}_x)(\forall F \in \mathcal{F})(F \cap N \neq \emptyset)\} \tag{A.24}$$

are respectively the sets of limit points and adherent points of  $\mathcal{F}$ <sup>28</sup>

A comparison of Eqs. (A.12) and (A.16) with Eqs. (A.23) and (A.24) respectively demonstrate their formal similarity; this inter-relation between filters and nets will be made precise in Definitions A.1.10 and A.1.11 below. It should be clear from the preceding two equations that

$$\lim(\mathcal{F}) \subseteq \text{adh}(\mathcal{F}), \tag{A.26}$$

with a similar result

$$\lim(\chi) \subseteq \text{adh}(\chi) \tag{A.27}$$

holding for nets because of the duality between nets and filters as displayed by Definitions A.1.9 and A.1.10 below, with the equality in Eqs. (A.26) and (A.27) being true (but not characterizing) for ultrafilters and ultranets respectively, see Example 4.2(3) for an account of this notion. It should be clear from the equations of Definition A.1.8 that

$$\text{adh}(\mathcal{F}) = \{x \in X : (\exists \text{ a finer filter } \mathcal{G} \supseteq \mathcal{F} \text{ on } X)(\mathcal{G} \rightarrow x)\} \tag{A.28}$$

consists of all the points of  $X$  to which some finer filter  $\mathcal{G}$  (in the sense that  $\mathcal{F} \subseteq \mathcal{G}$  implies every element of  $\mathcal{F}$  is also in  $\mathcal{G}$ ) converges in  $X$ ; thus

$$\text{adh}(\mathcal{F}) = \bigcup \lim(\mathcal{G} : \mathcal{G} \supseteq \mathcal{F}),$$

which corresponds to the net-result of Theorem A.1.5 below, that a net  $\chi$  adheres to  $x$  iff there is some subnet of  $\chi$  that converges to  $x$  in  $X$ . Thus if  $\zeta \preceq \chi$  is a subnet of  $\chi$  and  $\mathcal{F} \subseteq \mathcal{G}$  is a filter coarser than  $\mathcal{G}$  then

$$\begin{aligned} \lim(\chi) &\subseteq \lim(\zeta) & \lim(\mathcal{F}) &\subseteq \lim(\mathcal{G}) \\ \text{adh}(\zeta) &\subseteq \text{adh}(\chi) & \text{adh}(\mathcal{G}) &\subseteq \text{adh}(\mathcal{F}); \end{aligned}$$

a filter  $\mathcal{G}$  finer than a given filter  $\mathcal{F}$  corresponds to a subnet  $\zeta$  of a given net  $\chi$ . The implication of this correspondence should be clear from the association between nets and filters contained in Definitions A.1.10 and A.1.11.

A filter-base in  $X$  is a *non-empty* family  $(B_\alpha)_{\alpha \in \mathbb{D}} = {}_F\mathcal{B}$  of subsets of  $X$  characterized by

(FB1) There are no empty sets in the collection  ${}_F\mathcal{B}$ :  $(\forall \alpha \in \mathbb{D})(B_\alpha \neq \emptyset)$

(FB2) The intersection of any two members of  ${}_F\mathcal{B}$  contains another member of  ${}_F\mathcal{B}$ :  $B_\alpha, B_\beta \in {}_F\mathcal{B} \Rightarrow (\exists B \in {}_F\mathcal{B} : B \subseteq B_\alpha \cap B_\beta)$ ;

hence any class of subsets of  $X$  that does not contain the empty set and is closed under finite intersections is a base for a unique filter on  $X$ ; compare the properties (NB1) and (NB2) of a local basis given at the beginning of this Appendix. Similar to Definition A.1.1 for the local base, it is possible to define

**Definition A.1.9.** A filter-base  ${}_F\mathcal{B}$  in a set  $X$  is a subcollection of the filter  $\mathcal{F}$  on  $X$  having the property that each  $F \in \mathcal{F}$  contains some member of  ${}_F\mathcal{B}$ . Thus

$${}_F\mathcal{B} \stackrel{\text{def}}{=} \{B \in \mathcal{F} : B \subseteq F \text{ for each } F \in \mathcal{F}\} \quad (\text{A.29})$$

determines the filter

$$\mathcal{F} = \{F \subseteq X : B \subseteq F \text{ for some } B \in {}_F\mathcal{B}\} \quad (\text{A.30})$$

reciprocally as all supersets of the basic elements.

This is the smallest filter on  $X$  that contains  ${}_F\mathcal{B}$  and is said to be *the filter generated by its filter-base  ${}_F\mathcal{B}$* ; alternatively  ${}_F\mathcal{B}$  is the filter-base of  $\mathcal{F}$ . The entire neighborhood system  $\mathcal{N}_x$ , the local base  $\mathcal{B}_x$ ,  $\mathcal{N}_x \cap A$  for  $x \in \text{Cl}(A)$ , and the set of all residuals of a directed set  $\mathbb{D}$  are among the most useful examples of filter-bases on  $X$ ,  $A$  and  $\mathbb{D}$  respectively. Of course, every filter is trivially a filter-base of itself, and *the singletons  $\{\{x\}\}$ ,  $\{A\}$  are filter-bases that generate the principal filters  ${}_F\mathcal{P}(x)$  and  ${}_F\mathcal{P}(A)$  at  $x$ , and  $A$  respectively.*

Paralleling the case of topological subbase  ${}_T\mathcal{S}$ , a filter subbase  ${}_F\mathcal{S}$  can be defined on  $X$  to be any collection of subsets of  $X$  *with the finite intersection property* (as compared with  ${}_T\mathcal{S}$  where no such condition was necessary, this represents the fundamental point of departure between topology and filter) and it is not difficult to deduce that the filter generated by  ${}_F\mathcal{S}$  on  $X$  is obtained by taking all finite intersections  ${}_F\mathcal{S}_\Delta$  of members of  ${}_F\mathcal{S}$  followed by their

supersets  ${}_F\mathcal{S}_{\Sigma\Delta}$ .  $\mathcal{F}({}_F\mathcal{S}) := {}_F\mathcal{S}_{\Sigma\Delta}$  is the smallest filter on  $X$  that contains  ${}_F\mathcal{S}$  and is the filter *generated by  ${}_F\mathcal{S}$* .

Equation (A.24) can be put in the more useful and transparent form given by

**Theorem A.1.3.** For a filter  $\mathcal{F}$  in a space  $(X, \mathcal{T})$

$$\begin{aligned} \text{adh}(\mathcal{F}) &= \bigcap_{F \in \mathcal{F}} \text{Cl}(F) \\ &= \bigcap_{B \in {}_F\mathcal{B}} \text{Cl}(B), \end{aligned} \quad (\text{A.31})$$

and dually  $\text{adh}(\chi)$ , are closed sets.

*Proof.* Follows immediately from the definitions for the closure of a set Eq. (20) and the adherence of a filter Eq. (A.24). As always, it is a matter of convenience in using the basic filters  ${}_F\mathcal{B}$  instead of  $\mathcal{F}$  to generate the adherence set. ■

It is in fact true that the limit sets  $\text{lim}(\mathcal{F})$  and  $\text{lim}(\chi)$  are also closed set of  $X$ ; the arguments involving ultrafilters are omitted.

Similar to the notion of the adherence set of a filter is its *core* — a concept that unlike the adherence, is purely set-theoretic being the infimum of the filter and is not linked with any topological structure of the underlying (infinite) set  $X$  — defined as

$$\text{core}(\mathcal{F}) = \bigcap_{F \in \mathcal{F}} F. \quad (\text{A.32})$$

From Theorem A.1.3 and the fact that the closure of a set  $A$  is the smallest closed set that contains  $A$ , see Eq. (25) at the end of Tutorial 4, it is clear that in terms of filters

$$\begin{aligned} A &= \text{core}({}_F\mathcal{P}(A)) \\ \text{Cl}(A) &= \text{adh}({}_F\mathcal{P}(A)) \\ &= \text{core}(\text{Cl}({}_F\mathcal{P}(A))) \end{aligned} \quad (\text{A.33})$$

where  ${}_F\mathcal{P}(A)$  is the principal filter at  $A$ ; thus *the core and adherence sets of the principal filter at  $A$  are equal respectively to  $A$  and  $\text{Cl}(A)$*  — a classic example of equality in the general relation  $\text{Cl}(\bigcap A_\alpha) \subseteq \bigcap \text{Cl}(A_\alpha)$  — but both are empty, for example, in the case of an infinitely decreasing family of rationals centered at any irrational (leading to a principal filter-base of rationals at the chosen irrational). This is an important example demonstrating that *the infinite intersection of a non-empty family of (closed) sets with the finite intersection property may be empty, a situation that cannot arise on a finite set or an infinite compact set.* Filters on

$X$  with an empty core are said to be *free*, and are *fixed* otherwise: notice that by its very definition filters cannot be free on a finite set, and a free filter represents an additional feature that may arise in passing from finite to infinite sets. Clearly  $(\text{adh}(\mathcal{F}) = \emptyset) \Rightarrow (\text{core}(\mathcal{F}) = \emptyset)$ , but as the important example of the rational space in the reals illustrate, the converse need not be true. Another example of a free filter of the same type is provided by the filter-base  $\{[a, \infty) : a \in \mathbb{R}\}$  in  $\mathbb{R}$ . Both these examples illustrate the important property that *a filter is free iff it contains the cofinite filter*, and the cofinite filter is the smallest possible free filter on an infinite set. The free cofinite filter, as these examples illustrate, may be typically generated as follows. Let  $A$  be a subset of  $X$ ,  $x \in \text{Bdy}_{X-A}(A)$ , and consider the directed set Eq. (A.21) to generate the corresponding net in  $A$  given by  $\chi(N \in \mathcal{N}_x, t) = t \in A$ . Quite clearly, the core of any Frechet filter based on this net must be empty as the point  $x$  does not lie in  $A$ . In general, the intersection is empty because if it were not so then the complement of the intersection — which is an element of the filter — would be infinite in contravention of the hypothesis that the filter is Frechet. It should be clear that every filter finer than a free filter is also free, and any filter coarser than a fixed filter is fixed.

Nets and filters are complimentary concepts and one may switch from one to the other as follows.

**Definition A.1.10.** Let  $\mathcal{F}$  be a filter on  $X$  and let  $\mathbb{D}F_x = \{(F, x) : (F \in \mathcal{F})(x \in F)\}$  be a directed set with its natural direction  $(F, x) \preceq (G, y) \Rightarrow (G \subseteq F)$ . The net  $\chi_{\mathcal{F}}: \mathbb{D}F_x \rightarrow X$  defined by

$$\chi_{\mathcal{F}}(F, x) = x$$

is said to be associated with the filter  $\mathcal{F}$ , see Eq. (A.20).

**Definition A.1.11.** Let  $\chi : \mathbb{D} \rightarrow X$  be a net and  $\mathbb{R}_\alpha = \{\beta \in \mathbb{D} : \beta \succeq \alpha \in \mathbb{D}\}$  a residual in  $\mathbb{D}$ . Then

$${}_F\mathcal{B}_\chi \stackrel{\text{def}}{=} \{\chi(\mathbb{R}_\alpha) : \text{Res}(\mathbb{D}) \rightarrow X \text{ for all } \alpha \in \mathbb{D}\}$$

is the filter-base associated with  $\chi$ , and the corresponding filter  $\mathcal{F}_\chi$  obtained by taking all supersets of the elements of  ${}_F\mathcal{B}_\chi$  is the filter associated with  $\chi$ .

${}_F\mathcal{B}_\chi$  is a filter-base in  $X$  because  $\chi(\bigcap \mathbb{R}_\alpha) \subseteq \bigcap \chi(\mathbb{R}_\alpha)$ , that holds for any functional relation, proves (FB2). It is not difficult to verify that

(i)  $\chi$  is eventually in  $A \Rightarrow A \in \mathcal{F}_\chi$ , and

(ii)  $\chi$  is frequently in  $A \Rightarrow (\forall \mathbb{R}_\alpha \in \text{Res}(\mathbb{D})) (A \cap \chi(\mathbb{R}_\alpha) \neq \emptyset) \Rightarrow A \cap \mathcal{F}_\chi \neq \emptyset$ .

Limits and adherences are obviously preserved in switching between nets (respectively, filters) and the filters (respectively, nets) that they generate:

$$\lim(\chi) = \lim(\mathcal{F}_\chi), \quad \text{adh}(\chi) = \text{adh}(\mathcal{F}_\chi) \quad (\text{A.34})$$

$$\lim(\mathcal{F}) = \lim(\chi_{\mathcal{F}}), \quad \text{adh}(\mathcal{F}) = \text{adh}(\chi_{\mathcal{F}}). \quad (\text{A.35})$$

The proofs of the two parts of Eq. (A.34), for example, go respectively as follows.  $x \in \lim(\chi) \Leftrightarrow \chi$  is eventually in  $\mathcal{N}_x \Leftrightarrow (\forall N \in \mathcal{N}_x)(\exists F \in \mathcal{F}_\chi)$  such that  $(F \subseteq N) \Leftrightarrow x \in \lim(\mathcal{F}_\chi)$ , and  $x \in \text{adh}(\chi) \Leftrightarrow \chi$  is frequently in  $\mathcal{N}_x \Leftrightarrow (\forall N \in \mathcal{N}_x)(\forall F \in \mathcal{F}_\chi)(N \cap F \neq \emptyset) \Leftrightarrow x \in \text{adh}(\mathcal{F}_\chi)$ ; here  $F$  is a superset of  $\chi(\mathbb{R}_\alpha)$ .

Some examples of convergence of filters are

- (1) Any filter on an indiscrete space  $X$  converges to every point of  $X$ .
- (2) Any filter on a space that coincides with its topology (minus the empty set, of course) converges to every point of the space.
- (3) For each  $x \in X$ , the neighborhood filter  $\mathcal{N}_x$  converges to  $x$ ; this is the smallest filter on  $X$  that converges to  $x$ .
- (4) The *indiscrete* filter  $\mathcal{F} = \{X\}$  converges to no point in the space  $(X, \{\emptyset, A, X - A, X\})$ , but converges to every point of  $X - A$  if  $X$  has the topology  $\{\emptyset, A, X\}$  because the only neighborhood of any point in  $X - A$  is  $X$  which is contained in the filter.

One of the most significant consequences of convergence theory of sequences and nets, as shown by the two theorems and the corollary following, is that this can be used to describe the topology of a set. The proofs of the theorems also illustrate the close inter-relationship between nets and filters.

**Theorem A.1.4.** For a subset  $A$  of a topological space  $X$ ,

$$\text{Cl}(A) = \{x \in X : (\exists \text{ a net } \chi \text{ in } A)(\chi \rightarrow x)\}. \quad (\text{A.36})$$

*Proof. Necessity.* For  $x \in \text{Cl}(A)$ , construct a net  $\chi \rightarrow x$  in  $A$  as follows. Let  $\mathcal{B}_x$  be a topological local base at  $x$ , which by definition is the collection of all open sets of  $X$  containing  $x$ . For each  $\beta \in \mathbb{D}$ , the sets

$$N_\beta = \bigcap_{\alpha \preceq \beta} \{B_\alpha : B_\alpha \in \mathcal{B}_x\}$$

form a nested decreasing local neighborhood filter base at  $x$ . With respect to the directed set  $\mathbb{D}N_\beta = \{(N_\beta, \beta) : (\beta \in \mathbb{D})(x_\beta \in N_\beta)\}$  of Eq. (A.10), define the desired net in  $A$  by

$$\chi(N_\beta, \beta) = x_\beta \in N_\beta \bigcap A$$

where the family of non-empty decreasing subsets  $N_\beta \bigcap A$  of  $X$  constitute the filter-base in  $A$  as required by the directed set  $\mathbb{D}N_\beta$ . It now follows from Eq. (A.11) and the arguments in Example A.1.3(3) that  $x_\beta \rightarrow x$ ; compare the directed set of Eq. (A.21) for a more compact, yet essentially identical, argument. Carefully observe the dual roles of  $\mathcal{N}_x$  as a neighborhood filter base at  $x$ .

*Sufficiency.* Let  $\chi$  be a net in  $A$  that converges to  $x \in X$ . For any  $N_\alpha \in \mathcal{N}_x$ , there is a  $\mathbb{R}_\alpha \in \text{Res}(\mathbb{D})$  of Eq. (A.13) such that  $\chi(\mathbb{R}_\alpha) \subseteq N_\alpha$ . Hence the point  $\chi(\alpha) = x_\alpha$  of  $A$  belongs to  $N_\alpha$  so that  $A \bigcap N_\alpha \neq \emptyset$  which means, from Eq. (20), that  $x \in \text{Cl}(A)$ . ■

**Corollary.** *Together with Eqs. (20) and (22), it follows that*

$$\text{Der}(A) = \{x \in X : (\exists \text{ a net } \zeta \text{ in } A - \{x\})(\zeta \rightarrow x)\} \tag{A.37}$$

*The filter forms of Eqs. (A.36) and (A.37)*

$$\begin{aligned} \text{Cl}(A) &= \{x \in X : (\exists \text{ a filter } \mathcal{F} \text{ on } X) \\ &\quad (A \in \mathcal{F})(\mathcal{F} \rightarrow x)\} \\ \text{Der}(A) &= \{x \in X : (\exists \text{ a filter } \mathcal{F} \text{ on } X) \\ &\quad (A - \{x\} \in \mathcal{F})(\mathcal{F} \rightarrow x)\} \end{aligned} \tag{A.38}$$

*then follows from Eq. (A.25) and the finite intersection property (F2) of  $\mathcal{F}$  so that every neighborhood of  $x$  must intersect  $A$  (respectively  $A - \{x\}$ ) in Eq. (A.38) to produce the converging net needed in the proof of Theorem A.1.3.*

We end this discussion of convergence in topological spaces with a proof of the following theorem which demonstrates the relationship that “eventually in” and “frequently in” bears with each other; Eq. (A.39) below is the net-counterpart of the filter equation (A.28).

**Theorem A.1.5.** *If  $\chi$  is a net in a topological space  $X$ , then  $x \in \text{adh}(\chi)$  iff some subnet  $\zeta(\beta) = \chi(\sigma(\beta))$  of  $\chi(\alpha)$ , with  $\alpha \in \mathbb{D}$  and  $\beta \in \mathbb{E}$ , converges in  $X$  to  $x$ ; thus*

$$\text{adh}(\chi) = \{x \in X : (\exists \text{ a subnet } \zeta \preceq \chi \text{ in } X)(\zeta \rightarrow x)\}. \tag{A.39}$$

*Proof. Necessity.* Let  $x \in \text{adh}(\chi)$ . Define a subnet function  $\sigma : \mathbb{D}N_\alpha \rightarrow \mathbb{D}$  by  $\sigma(N_\alpha, \alpha) = \alpha$  where  $\mathbb{D}N_\alpha$  is the directed set of Eq. (A.10): (SN1) and (SN2) are quite evidently satisfied according to Eq. (A.11). Proceeding as in the proof of the preceding theorem it follows that  $x_\beta = \chi(\sigma(N_\alpha, \alpha)) = \zeta(N_\alpha, \alpha) \rightarrow x$  is the required converging subnet that exists from Eq. (A.15) and the fact that  $\chi(\mathbb{R}_\alpha) \bigcap N_\alpha \neq \emptyset$  for every  $N_\alpha \in \mathcal{N}_x$ , by hypothesis.

*Sufficiency.* Assume now that  $\chi$  has a subnet  $\zeta(N_\alpha, \alpha)$  that converges to  $x$ . If  $\chi$  does not adhere at  $x$ , there is a neighborhood  $N_\alpha$  of  $x$  not frequented by it, in which case  $\chi$  must be eventually in  $X - N_\alpha$ . Then  $\zeta(N_\alpha, \alpha)$  is also eventually in  $X - N_\alpha$  so that  $\zeta$  cannot be eventually in  $N_\alpha$ , a contradiction of the hypothesis that  $\zeta(N_\alpha, \alpha) \rightarrow x$ .<sup>29</sup> ■

Equations (A.36) and (A.39) imply that the closure of a subset  $A$  of  $X$  is the class of  $X$ -adherences of all the (sub)nets of  $X$  that are eventually in  $A$ . This includes both the constant nets yielding the isolated points of  $A$  and the non-constant nets leading to the cluster points of  $A$ , and implies the following physically useful relationship between convergence and topology that can be used as defining criteria for open and closed sets having a more appealing physical significance than the original definitions of these terms. Clearly, the term “net” is justifiably used here to include the subnets too.

The following corollary of Theorem A.1.5 summarizes the basic topological properties of sets in terms of nets (respectively, filters).

**Corollary.** *Let  $A$  be a subset of a topological space  $X$ . Then*

- (1)  *$A$  is closed in  $X$  iff every convergent net of  $X$  that is eventually in  $A$  actually converges to a point in  $A$  (respectively, iff the adhering*

<sup>29</sup>In a first countable space, while the corresponding proof of the first part of the theorem for sequences is essentially the same as in the present case, the more direct proof of the converse illustrates how the convenience of nets and directed sets may require more general arguments. Thus if a sequence  $(x_i)_{i \in \mathbb{N}}$  has a subsequence  $(x_{i_k})_{k \in \mathbb{N}}$  converging to  $x$ , then a more direct line of reasoning proceeds as follows. Since the subsequence converges to  $x$ , its tail  $(x_{i_k})_{k \geq j}$  must be in every neighborhood  $N$  of  $x$ . But as the number of such terms is infinite whereas  $\{i_k : k < j\}$  is only finite, it is necessary that for any given  $n \in \mathbb{N}$ , cofinitely many elements of the sequence  $(x_{i_k})_{i_k \geq n}$  be in  $N$ . Hence  $x \in \text{adh}((x_i)_{i \in \mathbb{N}})$ .

points of each filter-base on  $A$  all belong to  $A$ ). Thus no  $X$ -convergent net in a closed subset may converge to a point outside it.

- (2)  $A$  is open in  $X$  iff every convergent net of  $X$  that converges to a point in  $A$  is eventually in  $A$ . Thus no  $X$ -convergent net outside an open subset may converge to a point in the set.
- (3)  $A$  is closed-and-open (clopen) in  $X$  iff every convergent net of  $X$  that converges in  $A$  is eventually in  $A$  and conversely.
- (4)  $x \in \text{Der}(A)$  iff some net (respectively, filter-base) in  $A - \{x\}$  converges to  $x$ ; this clearly eliminates the isolated points of  $A$  and  $x \in \text{Cl}(A)$  iff some net (respectively, filter-base) in  $A$  converges to  $x$ .

*Remark.* The differences in these characterizations should be fully appreciated: If we consider the cluster points  $\text{Der}(A)$  of a net  $\chi$  in  $A$  as the *resource generated by*  $\chi$ , then a closed subset of  $X$  can be considered to be *selfish* as it keeps all its resource to itself:  $\text{Der}(A) \cap A = \text{Der}(A)$ . The opposite of this is a *donor* set that donates all its generated resources to its neighbor:  $\text{Der}(A) \cap X - A = \text{Der}(A)$ , while for a *neutral* set, both  $\text{Der}(A) \cap A \neq \emptyset$  and  $\text{Der}(A) \cap X - A \neq \emptyset$  implying that the convergence resources generated in  $A$  and  $X - A$  can be deposited only in the respective sets. The clopen sets (see diagram 2-2 of Fig. 22) are of some special interest as they are boundary less so that no net-resources can be generated in this case as any such limit are required to be simultaneously in the set and its complement.

**Example A.1.2.** (Continued). This continuation Example A.1.2 illustrates how sequential convergence is inadequate in spaces that are not first countable like the uncountable set with cocountable topology. In this topology, a sequence can converge to a point  $x$  in the space iff it has only a finite number of distinct terms, and is therefore eventually constant. Indeed, let the complement

$$G \stackrel{\text{def}}{=} X - F, \quad F = \{x_i : x_i \neq x, \quad i \in \mathbb{N}\}$$

of the countably closed sequential set  $F$  be an open neighborhood of  $x \in X$ . Because a sequence  $(x_i)_{i \in \mathbb{N}}$  in  $X$  converges to a point  $x \in X$  iff it is eventually in every neighborhood (including  $G$ ) of  $x$ , the sequence represented by the set  $F$  cannot converge

to  $x$  unless it is of the uncountable type<sup>30</sup>

$$(x_0, x_1, \dots, x_I, x_{I+1}, x_{I+1}, \dots) \tag{A.40}$$

with only a finite number  $I$  of distinct terms actually belonging to the closed sequential set  $F = X - G$ , and  $x_{I+1} = x$ . Note that as we are concerned only with the eventual behavior of the sequence, we may discard all distinct terms from  $G$  by considering them to be in  $F$ , and retain only the constant sequence  $(x, x, \dots)$  in  $G$ . In comparison with the cofinite case that was considered in Sec. 4, the entire countably infinite sequence can now lie outside a neighborhood of  $x$  thereby enforcing the eventual constancy of the sequence. This leads to a generalization of our earlier cofinite result in the sense that a cocountable filter on a cocountable space converges to every point in the space.

It is now straightforward to verify that for a point  $x_0$  in an uncountable cocountable space  $X$

- (a) Even though no sequence in the open set  $G = X - \{x_0\}$  can converge to  $x_0$ , yet  $x_0 \in \text{Cl}(G)$  since the intersection of any (uncountable) open neighborhood  $U$  of  $x_0$  with  $G$ , being an uncountable set, is not empty.
- (b) By Corollary 1 of Theorem A.1.5, the uncountable open set  $G = X - \{x_0\}$  is also closed in  $X$  because if any sequence  $(x_1, x_2, \dots)$  in  $G$  converges to some  $x \in X$ , then  $x$  must be in  $G$  as the sequence must be eventually constant in order for it to converge. But this is a contradiction as  $G$  cannot be closed since it is not countable.<sup>31</sup> By the same reckoning, although  $\{x_0\}$  is not an open set because its complement is not countable, nevertheless it follows from Eq. (A.40) that should any sequence converge to the only point  $x_0$  of this set, then it must eventually be in  $\{x_0\}$  so by Corollary 2 of the same theorem,  $\{x_0\}$  becomes an open set.
- (c) The identity map  $\mathbf{1} : X \rightarrow X_d$ , where  $X_d$  is  $X$  with discrete topology, is not continuous because the inverse image of any singleton of  $X_d$  is not open in  $X$ . Yet if a sequence converges in  $X$  to  $x$ , then its image  $(\mathbf{1}(x)) = (x)$  must actually converge to  $x$  in  $X_d$  because a sequence converges in a discrete space, as in the cofinite or cocountable spaces, iff it is eventually constant; this is so because each element of a discrete space being clopen is boundary-less.

<sup>30</sup>This is uncountable because interchanging any two eventual terms of the sequence does not alter the sequence.

<sup>31</sup>Note that  $\{x\}$  is a 1-point set but  $(x)$  is an uncountable sequence.

This pathological behavior of sequences in a non Hausdorff, non first countable space does not arise if the discrete indexing set of sequences is replaced by a continuous, uncountable directed set like  $\mathbb{R}$  for example, leading to nets in place of sequences. In this case the net can be in an open set without having to be constant valued in order to converge to a point in it as the open set can be defined as the complement of a closed countable part of the uncountable net. The careful reader could not have failed to notice that the burden of the above arguments, as also of that in the example following Theorem 4.6, is to formalize the fact that since a closed set is already defined as a countable (respectively finite) set, the closure operation cannot add further points to it from its complement, and any sequence that converges in an open set in these topologies must necessarily be eventually constant at its point of convergence, a restriction that no longer applies to a net. The cocountable topology thus has the very interesting property of filtering out a countable part from an uncountable set, as for example the rationals in  $\mathbb{R}$ .

This example serves to illustrate the hard truth that in a space that is not first countable, the simplicity of sequences is not enough to describe its topological character, and in fact “sequential convergence will be able to describe only those topologies in which the number of (basic) neighborhoods around each point is no greater than the number of terms in the sequences”, [Willard, 1970]. It is important to appreciate the significance of this interplay of convergence of sequences and nets (and of continuity of functions of Appendix A.1) and the topology of the underlying spaces.

A comparison of the defining properties (T1)–(T3) of topology  $\mathcal{T}$  with (F1)–(F3) of that of the filter  $\mathcal{F}$ , shows that a filter is very close to a topology with the main difference being with regard to the empty set which must always be in  $\mathcal{T}$  but never in  $\mathcal{F}$ . Addition of the empty set to a filter yields a topology, but removal of the empty set from a topology need not produce the corresponding filter as the topology may contain nonintersecting sets.

The distinction between the topological and filter-bases should be carefully noted. Thus

- (a) While the topological base may contain the empty set, a filter-base cannot.
- (b) From a given topology, form a common base by

dropping all basic open sets that do not intersect. Then a (coarser) topology can be generated from this base by taking all unions, and a filter by taking all supersets according to Eq. (A.30). For any given filter this expression may be used to extract a subclass  ${}_F\mathcal{B}$  as a base for  $\mathcal{F}$ .

## A.2. Initial and Final Topology

The commutative diagram of Fig. contains four sub-diagrams  $X - X_B - f(X)$ ,  $Y - X_B - f(X)$ ,  $X - X_B - Y$  and  $X - f(X) - Y$ . Of these, the first two are especially significant as they can be used to conveniently define the topologies on  $X_B$  and  $f(X)$  from those of  $X$  and  $Y$ , so that  $f_B$ ,  $f_B^{-1}$  and  $G$  have some desirable continuity properties; we recall that a function  $f : X \rightarrow Y$  is continuous if inverse images of open sets of  $Y$  are open in  $X$ . This simple notion of continuity needs refinement in order that topologies on  $X_B$  and  $f(X)$  be unambiguously defined from those of  $X$  and  $Y$ , a requirement that leads to the concepts of the so-called *final* and *initial topologies*. To appreciate the significance of these new constructs, note that if  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  is a continuous function, there may be open sets in  $X$  that are not inverse images of open — or for that matter of any — subset of  $Y$ , just as it is possible for non-open subsets of  $Y$  to contribute to  $\mathcal{U}$ . When the triple  $\{\mathcal{U}, f, \mathcal{V}\}$  are tuned in such a manner that these are impossible, the topologies so generated on  $X$  and  $Y$  are the initial and final topologies respectively; they are the smallest (coarsest) and largest (finest) topologies on  $X$  and  $Y$  that make  $f : X \rightarrow Y$  continuous. It should be clear that every image and preimage continuous function is continuous, but the converse is not true.

Let  $\text{sat}(U) := f^{-1}f(U) \subseteq X$  be the saturation of an open set  $U$  of  $X$  and  $\text{comp}(V) := ff^{-1}(V) = V \cap f(X) \in Y$  be the component of an open set  $V$  of  $Y$  on the range  $f(X)$  of  $f$ . Let  $\mathcal{U}_{\text{sat}}$ ,  $\mathcal{V}_{\text{comp}}$  denote respectively the saturations  $\mathcal{U}_{\text{sat}} = \{\text{sat}(U) : U \in \mathcal{U}\}$  of the open sets of  $X$  and the components  $\mathcal{V}_{\text{comp}} = \{\text{comp}(V) : V \in \mathcal{V}\}$  of the open sets of  $Y$  whenever these are also open in  $X$  and  $Y$  respectively. Plainly,  $\mathcal{U}_{\text{sat}} \subseteq \mathcal{U}$  and  $\mathcal{V}_{\text{comp}} \subseteq \mathcal{V}$ .

**Definition A.2.1.** For a function  $e : X \rightarrow (Y, \mathcal{V})$ , the preimage or initial topology of  $X$  based on (generated by)  $e$  and  $\mathcal{V}$  is

$$\text{IT}\{e; \mathcal{V}\} \stackrel{\text{def}}{=} \{U \subseteq X : U = e^{-1}(V) \text{ if } V \in \mathcal{V}_{\text{comp}}\}, \quad (\text{A.41})$$



while for  $q : (X, \mathcal{U}) \rightarrow Y$ , the image or final topology of  $Y$  based on (generated by)  $\mathcal{U}$  and  $q$  is

$$\text{FT}\{\mathcal{U}; q\} \stackrel{\text{def}}{=} \{V \subseteq Y : q^{-1}(V) = U \text{ if } U \in \mathcal{U}_{\text{sat}}\}. \tag{A.42}$$

Thus, the topology of  $(X, \text{IT}\{e; \mathcal{V}\})$  consists of, and only of, the  $e$ -saturations of all the open sets of  $e(X)$ , while the open sets of  $(Y, \text{FT}\{\mathcal{U}; q\})$  are the  $q$ -images in  $Y$  (and not just in  $q(X)$ ) of all the  $q$ -saturated open sets of  $X$ .<sup>32</sup> The need for defining (A.41) in terms of  $\mathcal{V}_{\text{comp}}$  rather than  $\mathcal{V}$  will become clear in the following. The subspace topology  $\text{IT}\{i; \mathcal{U}\}$  of a subset  $A \subseteq (X, \mathcal{U})$  is a basic example of the initial topology by the inclusion map  $i : X \supseteq A \rightarrow (X, \mathcal{U})$ , and we take its generalization  $e : (A, \text{IT}\{e; \mathcal{V}\}) \rightarrow (Y, \mathcal{V})$  that embeds a subset  $A$  of  $X$  into  $Y$  as the prototype of a preimage continuous map. Clearly the topology of  $Y$  may also contain open sets not in  $e(X)$ , and any subset in  $Y - e(X)$  may be added to the topology of  $Y$  without altering the preimage topology of  $X$ : *open sets of  $Y$  not in  $e(X)$  may be neglected in obtaining the preimage topology* as  $e^{-1}(Y - e(X)) = \emptyset$ . The final topology on a quotient set by the quotient map  $Q : (X, \mathcal{U}) \rightarrow X/\sim$ , which is just the collection of  $Q$ -images of the  $Q$ -saturated open sets of  $X$ , known as the *quotient topology of  $X/\sim$* , is the basic example of the image topology and the resulting space  $(X/\sim, \text{FT}\{\mathcal{U}; Q\})$  is called the *quotient space*. We take the generalization  $q : (X, \mathcal{U}) \rightarrow (Y, \text{FT}\{\mathcal{U}; q\})$  of  $Q$  as the prototype of an image continuous function.

The following results are specifically useful in dealing with initial and final topologies; compare the corresponding results for open maps given later.

**Theorem A.2.1.** *Let  $(X, \mathcal{U})$  and  $(Y_1, \mathcal{V}_1)$  be topological spaces and let  $X_1$  be a set. If  $f : X_1 \rightarrow (Y_1, \mathcal{V}_1)$ ,  $q : (X, \mathcal{U}) \rightarrow X_1$ , and  $h = f \circ q : (X, \mathcal{U}) \rightarrow (Y_1, \mathcal{V}_1)$  are functions with the topology  $\mathcal{U}_1$  of  $X_1$  given by  $\text{FT}\{\mathcal{U}; q\}$ , then*

- (a)  $f$  is continuous iff  $h$  is continuous,
- (b)  $f$  is image continuous iff  $\mathcal{V}_1 = \text{FT}\{\mathcal{U}; h\}$ .

**Theorem A.2.2.** *Let  $(Y, \mathcal{V})$  and  $(X_1, \mathcal{U}_1)$  be topological spaces and let  $Y_1$  be a set. If  $f : (X_1, \mathcal{U}_1) \rightarrow Y_1$ ,  $e : Y_1 \rightarrow (Y, \mathcal{V})$  and  $g = e \circ f : (X_1, \mathcal{U}_1) \rightarrow (Y, \mathcal{V})$  are function with the topology  $\mathcal{V}_1$  of  $Y_1$  given by  $\text{IT}\{e; \mathcal{V}\}$ , then*

- (a)  $f$  is continuous iff  $g$  is continuous,
- (b)  $f$  is preimage continuous iff  $\mathcal{U}_1 = \text{IT}\{g; \mathcal{V}\}$ .

As we need the second part of these theorems in our applications, their proofs are indicated below. The special significance of the first parts is that they ensure the converse of the usual result that the composition of two continuous functions is continuous, namely that one of the components of a composition is continuous whenever the composition is so.

*Proof of Theorem A.2.1.* If  $f$  be image continuous,  $\mathcal{V}_1 = \{V_1 \subseteq Y_1 : f^{-1}(V_1) \in \mathcal{U}_1\}$  and  $\mathcal{U}_1 = \{U_1 \subseteq X_1 : q^{-1}(U_1) \in \mathcal{U}\}$  are the final topologies of  $Y_1$  and  $X_1$  based on the topologies of  $X_1$  and  $X$ , respectively. Then  $\mathcal{V}_1 = \{V_1 \subseteq Y_1 : q^{-1}f^{-1}(V_1) \in \mathcal{U}\}$  shows that  $h$  is image continuous.

Conversely, when  $h$  is image continuous,  $\mathcal{V}_1 = \{V_1 \subseteq Y_1 : h^{-1}(V_1) \in \mathcal{U}\} = \{V_1 \subseteq Y_1 : q^{-1}f^{-1}(V_1) \in \mathcal{U}\}$ , with  $\mathcal{U}_1 = \{U_1 \subseteq X_1 : q^{-1}(U_1) \in \mathcal{U}\}$ , proves  $f^{-1}(V_1)$  to be open in  $X_1$  and thereby  $f$  to be image continuous. ■

*Proof of Theorem A.2.2.* If  $f$  be preimage continuous,  $\mathcal{V}_1 = \{V_1 \subseteq Y_1 : V_1 = e^{-1}(V) \text{ if } V \in \mathcal{V}\}$  and  $\mathcal{U}_1 = \{U_1 \subseteq X_1 : U_1 = f^{-1}(V_1) \text{ if } V_1 \in \mathcal{V}_1\}$  are the initial topologies of  $Y_1$  and  $X_1$  respectively. Hence from  $\mathcal{U}_1 = \{U_1 \subseteq X_1 : U_1 = f^{-1}e^{-1}(V) \text{ if } V \in \mathcal{V}\}$  it follows that  $g$  is preimage continuous.

Conversely, when  $g$  is preimage continuous,  $\mathcal{U}_1 = \{U_1 \subseteq X_1 : U_1 = g^{-1}(V) \text{ if } V \in \mathcal{V}\} = \{U_1 \subseteq X_1 : U_1 = f^{-1}e^{-1}(V) \text{ if } V \in \mathcal{V}\}$  and  $\mathcal{V}_1 = \{V_1 \subseteq Y_1 : V_1 = e^{-1}(V) \text{ if } V \in \mathcal{V}\}$  show that  $f$  is preimage continuous. ■

Since both Eqs. (A.41) and (A.42) are in terms of inverse images (the first of which constitutes a direct, and the second an inverse, problem) the image  $f(U) = \text{comp}(V)$  for  $V \in \mathcal{V}$  is of interest as it indicates the relationship of the openness of  $f$  with its continuity. This, and other related concepts are examined below, where the range space  $f(X)$  is always taken to be a subspace of  $Y$ . Openness of a function  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  is the “inverse” of continuity, when images of open sets of  $X$  are required to be open in  $Y$ ; such a function is said to be *open*. Following are two of the important properties of open functions.

<sup>32</sup>We adopt the convention of denoting arbitrary preimage and image continuous functions by  $e$  and  $q$  respectively even though they are not injective or surjective; recall that the embedding  $e : X \supseteq A \rightarrow Y$  and the association  $q : X \rightarrow f(X)$  are 1 : 1 and onto respectively.

- (1) If  $f : (X, \mathcal{U}) \rightarrow (Y, f(\mathcal{U}))$  is an open function, then so is  $f_{<} : (X, \mathcal{U}) \rightarrow (f(X), \text{IT}\{i; f(\mathcal{U})\})$ . The converse is true if  $f(X)$  is an open set of  $Y$ ; thus openness of  $f_{<} : (X, \mathcal{U}) \rightarrow (f(X), f_{<}(\mathcal{U}))$  implies that of  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  whenever  $f(X)$  is open in  $Y$  such that  $f_{<}(U) \in \mathcal{V}$  for  $U \in \mathcal{U}$ . The truth of this last assertion follows easily from the fact that if  $f_{<}(U)$  is an open set of  $f(X) \subset Y$ , then necessarily  $f_{<}(U) = V \cap f(X)$  for some  $V \in \mathcal{V}$ , and the intersection of two open sets of  $Y$  is again an open set of  $Y$ .
- (2) If  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  and  $g : (Y, \mathcal{V}) \rightarrow (Z, \mathcal{W})$  are open functions then  $g \circ f : (X, \mathcal{U}) \rightarrow (Z, \mathcal{W})$  is also open. It follows that the condition in (1) on  $f(X)$  can be replaced by the requirement that the inclusion  $i : (f(X), \text{IT}\{i; \mathcal{V}\}) \rightarrow (Y, \mathcal{V})$  be an open map. This interchange of  $f(X)$  with its inclusion  $i : f(X) \rightarrow Y$  into  $Y$  is a basic result that finds application in many situations.

Collected below are some useful properties of the initial and final topologies that we need in this work.

**Initial Topology.** In Fig. 21(b), consider  $Y_1 = h(X_1)$ ,  $e \rightarrow i$  and  $f \rightarrow h_{<} : X_1 \rightarrow (h(X_1), \text{IT}\{i; \mathcal{V}\})$ . From  $h^-(B) = h^-(B \cap h(X_1))$  for any  $B \subseteq Y$ , it follows that for an open set  $V$  of  $Y$ ,  $h^-(V_{\text{comp}}) = h^-(V)$  is an open set of  $X_1$  which, if the topology of  $X_1$  is  $\text{IT}\{h; \mathcal{V}\}$ , are the only open sets of  $X_1$ . Because  $V_{\text{comp}}$  is an open set of  $h(X_1)$  in its subspace topology, this implies that the preimage topologies  $\text{IT}\{h; \mathcal{V}\}$  and  $\text{IT}\{h_{<}; \text{IT}\{i; \mathcal{V}\}\}$  of  $X_1$  generated by  $h$  and  $h_{<}$  are the same. Thus the preimage topology of  $X_1$  is not affected if  $Y$  is replaced by the subspace  $h(X_1)$ , the part  $Y - h(X_1)$  contributing nothing to  $\text{IT}\{h; \mathcal{V}\}$ .

A preimage continuous function  $e : X \rightarrow (Y, \mathcal{V})$  is not necessarily an open function. Indeed, if  $U = e^-(V) \in \text{IT}\{e; \mathcal{V}\}$ , it is almost trivial to verify along the lines of the restriction of open maps to its range, that  $e(U) = ee^-(V) = e(X) \cap V$ ,  $V \in \mathcal{V}$ , is open in  $Y$  (implying that  $e$  is an open map) iff  $e(X)$  is an open subset of  $Y$  (because finite intersections of open sets are open). A special case of this is the important consequence that the restriction  $e_{<} : (X, \text{IT}\{e; \mathcal{V}\}) \rightarrow (e(X), \text{IT}\{i; \mathcal{V}\})$  of  $e : (X, \text{IT}\{h; \mathcal{V}\}) \rightarrow (Y, \mathcal{V})$  to its range is an open map. Even though a preimage continuous map need not be open, it is true that an injective, continuous and open map  $f : X \rightarrow (Y, \mathcal{V})$  is preimage

continuous. Indeed, from its injectivity and continuity, inverse images of all open subsets of  $Y$  are saturated-open in  $X$ , and openness of  $f$  ensures that these are the only open sets of  $X$  the condition of injectivity being required to exclude non-saturated sets from the preimage topology. It is therefore possible to rewrite Eq. (A.41) as

$$U \in \text{IT}\{e; \mathcal{V}\} \Leftrightarrow e(U) = V \text{ if } V \in \mathcal{V}_{\text{comp}}, \quad (\text{A.43})$$

and to compare it with the following criterion for an injective, open-continuous map  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  that necessarily satisfies  $\text{sat}(A) = A$  for all  $A \subseteq X$

$$U \in \mathcal{U} \Leftrightarrow (\{f(U)\}_{U \in \mathcal{U}} = \mathcal{V}_{\text{comp}}) \wedge (f^{-1}(V)|_{V \in \mathcal{V}} \in \mathcal{U}). \quad (\text{A.44})$$

**Final Topology.** Since it is necessarily produced on the range  $\mathcal{R}(q)$  of  $q$ , the final topology is often considered in terms of a surjection. This however is not necessary as, much in the spirit of the initial topology,  $Y - q(X) \neq \emptyset$  inherits the discrete topology without altering anything, thereby allowing condition (A.42) to be restated in the following more transparent form

$$V \in \text{FT}\{\mathcal{U}; q\} \Leftrightarrow V = q(U) \text{ if } U \in \mathcal{U}_{\text{sat}}, \quad (\text{A.45})$$

and to compare it with the following criterion for a surjective, open-continuous map  $f : (X, \mathcal{U}) \rightarrow (Y, \mathcal{V})$  that necessarily satisfies  $fB = B$  for all  $B \subseteq Y$

$$V \in \mathcal{V} \Leftrightarrow (\mathcal{U}_{\text{sat}} = \{f^-(V)\}_{V \in \mathcal{V}}) \wedge (f(U)|_{U \in \mathcal{U}} \in \mathcal{V}). \quad (\text{A.46})$$

As may be anticipated from Fig. 21, the final topology does not behave as well for subspaces as the initial topology does. This is so because in Fig. 21(a) the two image continuous functions  $h$  and  $q$  are connected by a preimage continuous inclusion  $f$ , whereas in Fig. 21(b) all the three functions are preimage continuous. Thus quite like open functions, although image continuity of  $h : (X, \mathcal{U}) \rightarrow (Y_1, \text{FT}\{\mathcal{U}; h\})$  implies that of  $h_{<} : (X, \mathcal{U}) \rightarrow (h(X), \text{IT}\{i; \text{FT}\{\mathcal{U}; h\}\})$  for a subspace  $h(X)$  of  $Y_1$ , the converse need not be true unless — entirely like open functions again — either  $h(X)$  is an open set of  $Y_1$  or  $i : (h(X), \text{IT}\{i; \text{FT}\{\mathcal{U}; h\}\}) \rightarrow (X, \text{FT}\{\mathcal{U}; h\})$  is an open map. Since an open

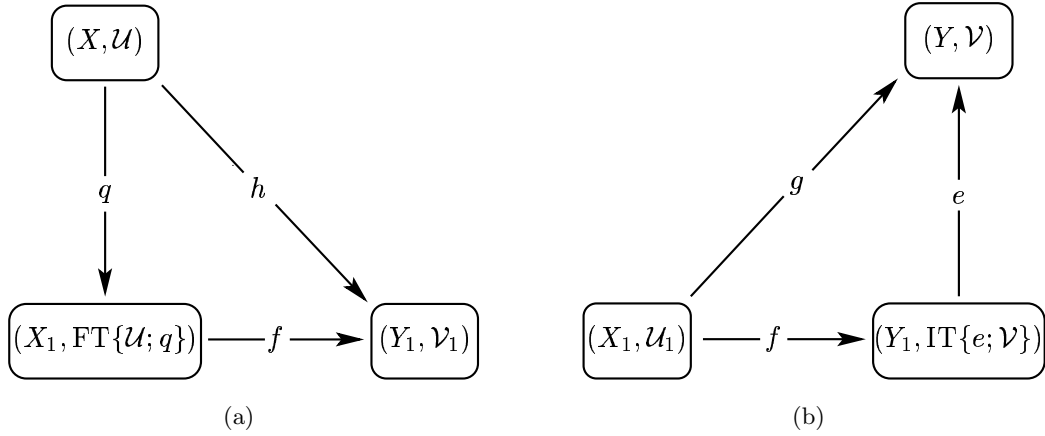


Fig. 21. Continuity in final and initial topologies.

preimage continuous map is image continuous, this makes  $i : h(X) \rightarrow Y_1$  an inital function and hence all the three legs of the commutative diagram image continuous.

Like preimage continuity, an image continuous function  $q : (X, \mathcal{U}) \rightarrow Y$  need not be open. However, although the restriction of an image continuous function to the saturated open sets of its domain is an open function,  $q$  is unrestrictedly open iff the saturation of every open set of  $X$  is also open in  $X$ . In fact it can be verified without much effort that a continuous, open surjection is image continuous.

Combining Eqs. (A.43) and (A.45) gives the following criterion for ininality

$$\begin{aligned}
 &U \text{ and } V \in \text{IFT}\{\mathcal{U}_{\text{sat}}; f; \mathcal{V}\} \\
 &\Leftrightarrow (\{f(U)\}_{U \in \mathcal{U}_{\text{sat}}} = \mathcal{V})(\mathcal{U}_{\text{sat}} = \{f^{-1}(V)\}_{V \in \mathcal{V}}),
 \end{aligned}
 \tag{A.47}$$

which reduces to the following for a homeomorphism  $f$  that satisfies both  $\text{sat}(A) = A$  for  $A \subseteq X$  and  $fB = B$  for  $B \subseteq Y$

$$\begin{aligned}
 &U \text{ and } V \in \text{HOM}\{\mathcal{U}; f; \mathcal{V}\} \\
 &\Leftrightarrow (\mathcal{U} = \{f^{-1}(V)\}_{V \in \mathcal{V}})(\{f(U)\}_{U \in \mathcal{U}} = \mathcal{V})
 \end{aligned}
 \tag{A.48}$$

and compares with

$$\begin{aligned}
 &U \text{ and } V \in \text{OC}\{\mathcal{U}; f; \mathcal{V}\} \\
 &\Leftrightarrow (\text{sat}(U) \in \mathcal{U} : \{f(U)\}_{U \in \mathcal{U}} = \mathcal{V}_{\text{comp}}) \\
 &\wedge (\text{comp}(V) \in \mathcal{V} : \{f^{-1}(V)\}_{V \in \mathcal{V}} = \mathcal{U}_{\text{sat}})
 \end{aligned}
 \tag{A.49}$$

for an open-continuous  $f$ .

The following is a slightly more general form of the restriction on the inclusion that is needed for image continuity to behave well for subspaces of  $Y$ .

**Theorem A.2.3.** *Let  $q : (X, \mathcal{U}) \rightarrow (Y, \text{FT}\{\mathcal{U}; q\})$  be an image continuous function. For a subspace  $B$  of  $(Y, \text{FT}\{\mathcal{U}; q\})$ ,*

$$\text{FT}\{\text{IT}\{j; \mathcal{U}\}; q_{<}\} = \text{IT}\{i; \text{FT}\{\mathcal{U}; q\}\}$$

where  $q_{<} : (q^{-1}(B), \text{IT}\{j; \mathcal{U}\}) \rightarrow (B, \text{FT}\{\text{IT}\{j; \mathcal{U}\}; q_{<}\})$ , if either  $q$  is an open map or  $B$  is an open set of  $Y$ .

In summary we have the useful result that an open preimage continuous function is image continuous and an open image continuous function is preimage continuous, where the second assertion follows on neglecting non-saturated open sets in  $X$ ; this is permitted in as far as the generation of the final topology is concerned, as these sets produce the same images as their saturations. Hence an image continuous function  $q : X \rightarrow Y$  is preimage continuous iff every open set in  $X$  is saturated with respect to  $q$ , and a preimage continuous function  $e : X \rightarrow Y$  is image continuous iff the  $e$ -image of every open set of  $X$  is open in  $Y$ .

### A.3. More on Topological Spaces

This Appendix — which completes the review of those concepts of topological spaces begun in Tutorial 4 that are needed for a proper understanding of this work — begins with the following summary of the different possibilities in the distribution of  $\text{Der}(A)$  and  $\text{Bdy}(A)$  between sets  $A \subseteq X$  and its complement  $X - A$ , and follows it up with a

few other important topological concepts that have been used, explicitly or otherwise, in this paper.

**Definition A.3.1** (Separation, Connected Space). A separation (disconnection) of  $X$  is a pair of mutually disjoint nonempty open (and therefore closed) subsets  $H_1$  and  $H_2$  such that  $X = H_1 \cup H_2$ . A space  $X$  is said to be connected if it has no separation, that is, if it cannot be partitioned into two open or two closed non-empty subsets.  $X$  is separated (disconnected) if it is not connected.

It follows from the definition, that for a disconnected space  $X$  the following are equivalent statements.

- There exist a pair of disjoint non-empty open subsets of  $X$  that cover  $X$ ,
- There exist a pair of disjoint non-empty closed subsets of  $X$  that cover  $X$ ,
- There exist a pair of disjoint non-empty clopen subsets of  $X$  that cover  $X$ ,
- There exists a non-empty, proper, clopen subset of  $X$ .

By a *connected subset* is meant a subset of  $X$  that is connected *when provided with its relative topology making it a subspace of  $X$* . Thus any connected subset of a topological space must necessarily be contained in any clopen set that might intersect it: if  $C$  and  $H$  are respectively connected and clopen subsets of  $X$  such that  $C \cap H \neq \emptyset$ , then  $C \subset H$  because  $C \cap H$  is a non-empty clopen set in  $C$  which must contain  $C$  because  $C$  is connected.

For testing whether a subset of a topological space is connected, the following relativized form of (a)–(d) is often useful.

**Lemma A.3.1.** *A subset  $A$  of  $X$  is disconnected iff there are disjoint open sets  $U$  and  $V$  of  $X$  satisfying*

$$\begin{aligned} U \cap A \neq \emptyset \neq V \cap A \text{ such that } A \subseteq U \cup V, \\ \text{with } U \cap V \cap A = \emptyset \end{aligned} \tag{A.50}$$

*or there are disjoint closed sets  $E$  and  $F$  of  $X$  satisfying*

$$\begin{aligned} E \cap A \neq \emptyset \neq F \cap A \text{ such that } A \subseteq E \cup F, \\ \text{with } E \cap F \cap A = \emptyset. \end{aligned} \tag{A.51}$$

*Thus  $A$  is disconnected iff there are disjoint clopen subsets in the relative topology of  $A$  that cover  $A$ .*

**Lemma A.3.2.** *If  $A$  is a subspace of  $X$ , a separation of  $A$  is a pair of disjoint nonempty subsets  $H_1$  and  $H_2$  of  $A$  whose union is  $A$  neither of which contains a cluster point of the other.  $A$  is connected iff there is no separation of  $A$ .*

*Proof.* Let  $H_1$  and  $H_2$  be a separation of  $A$  so that they are clopen subsets of  $A$  whose union is  $A$ . As  $H_1$  is a closed subset of  $A$  it follows that  $H_1 = \text{Cl}_X(H_1) \cap A$ , where  $\text{Cl}_X(H_1) \cap A$  is the closure of  $H_1$  in  $A$ ; hence  $\text{Cl}_X(H_1) \cap H_2 = \emptyset$ . But as the closure of a subset is the union of the set and its adherents, an empty intersection signifies that  $H_2$  cannot contain any of the cluster points of  $H_1$ . A similar argument shows that  $H_1$  does not contain any adherent of  $H_2$ .

Conversely suppose that neither  $H_1$  nor  $H_2$  contain an adherent of the other:  $\text{Cl}_X(H_1) \cap H_2 = \emptyset$  and  $\text{Cl}_X(H_2) \cap H_1 = \emptyset$ . Hence  $\text{Cl}_X(H_1) \cap A = H_1$  and  $\text{Cl}_X(H_2) \cap A = H_2$  so that both  $H_1$  and  $H_2$  are closed in  $A$ . But since  $H_1 = A - H_2$  and  $H_2 = A - H_1$ , they must also be open in the relative topology of  $A$ . ■

Following are some useful properties of connected spaces.

- The closure of any connected subspace of a space is connected. In general, every  $B$  satisfying

$$A \subseteq B \subseteq \text{Cl}(A)$$

is connected. Thus any subset of  $X$  formed from  $A$  by adjoining to it some or all of its adherents is connected so that a topological space with a dense connected subset is connected.

- The union of any class of connected subspaces of  $X$  with nonempty intersection is a connected subspace of  $X$ .
- A topological space is connected iff there is a covering of the space consisting of connected sets with nonempty intersection. Connectedness is a topological property: Any space homeomorphic to a connected space is itself connected.
- If  $H_1$  and  $H_2$  is a separation of  $X$  and  $A$  is any connected subset  $A$  of  $X$ , then either  $A \subseteq H_1$  or  $A \subseteq H_2$ .

While the real line  $\mathbb{R}$  is connected, a subspace of  $\mathbb{R}$  is connected iff it is an interval in  $\mathbb{R}$ .

		$X - A$		
		1. Donor	2. Selfish (Closed)	3. Neutral
$A$	1. Donor			
	2. Selfish (Closed)			
	3. Neutral			

Fig. 22. Classification of a subset  $A$  of  $X$  relative to the topology of  $X$ . The derived set of  $A$  may intersect both  $A$  and  $X - A$  (row 3), may be entirely in  $A$  (row 2), or may be wholly in  $X - A$  (row 1).  $A$  is closed iff  $\text{Bdy}(A) \subseteq A$  (row 2), open iff  $\text{Bdy}(A) \subseteq X - A$  (column 2), and clopen iff  $\text{Bdy}(A) = \emptyset$  when the derived sets of both  $A$  and  $X - A$  are contained in the respective sets. An open set, beside being closed, may also be neutral or donor.

The important concept of total disconnectedness introduced below needs the following

**Definition A.3.2** (Component). A component  $C^*$  of a space  $X$  is a maximally (with respect to inclusion) connected subset of  $X$ .

Thus a component is a connected subspace which is not properly contained in any larger connected subspace of  $X$ . The maximal element need not be unique as there can be more than one component of a given space and a “maximal” criterion rather than “maximum” is used as the component that need not contain every connected subset of  $X$ ; it sim-

ply must not be contained in any other connected subset of  $X$ . Components can be constructively defined as follows: Let  $x \in X$  be any point. Consider the collection of all connected subsets of  $X$  to which  $x$  belongs. Since  $\{x\}$  is one such a set, the collection is non-empty. As the intersection of the collection is non-empty, its union is a non-empty connected set  $C$ . This is the largest connected set containing  $x$  and is therefore a component containing  $x$  and we have

- (C1) Let  $x \in X$ . The unique component of  $X$  containing  $x$  is the union of all the connected subsets of  $X$  that contain  $x$ . Conversely any non-empty connected subset  $A$  of  $X$  is contained

in that unique component of  $X$  to which each of the points of  $A$  belong. Hence a topological space is connected iff it is the unique component of itself.

- (C2) Each component  $C^*$  of  $X$  is a closed set of  $X$ : By property (c1) above,  $\text{Cl}(C^*)$  is also connected and from  $C^* \subseteq \text{Cl}(C^*)$  it follows that  $C^* = \text{Cl}(C^*)$ . Components need not be open sets of  $X$ : an example of this is the space of rationals  $\mathbb{Q}$  in reals in which the components are the individual points which cannot be open in  $\mathbb{R}$ ; see Example 2 below.
- (C3) Components of  $X$  are equivalence classes of  $(X, \sim)$  with  $x \sim y$  iff they are in the same component: while reflexivity and symmetry are obvious enough, transitivity follows because if  $x, y \in C_1$  and  $y, z \in C_2$  with  $C_1, C_2$  connected subsets of  $X$ , then  $x$  and  $z$  are in the set  $C_1 \cup C_2$  which is connected by property c(2) above as they have the point  $y$  in common. Components are connected disjoint subsets of  $X$  whose union is  $X$  (i.e. they form a partition of  $X$  with each point of  $X$  contained in exactly one component of  $X$ ) such that any connected subset of  $X$  can be contained in only one of them. Because a connected subspace cannot contain in it any clopen subset of  $X$ , it follows that *every clopen connected subspace must be a component of  $X$* .

Even when a space is disconnected, it is always possible to decompose it into pairwise disjoint connected subsets. If  $X$  is a discrete space this is the only way in which  $X$  may be decomposed into connected pieces. If  $X$  is not discrete, there may be other ways of doing this. For example, the space

$$X = \{x \in \mathbb{R} : (0 \leq x \leq 1) \vee (2 < x < 3)\}$$

has the following distinct decomposition into three connected subsets:

$$X = \left[0, \frac{1}{2}\right) \cup \left[\frac{1}{2}, 1\right] \cup \left(2, \frac{7}{3}\right] \cup \left(\frac{7}{3}, 3\right)$$

$$X = \{0\} \cup \left(\bigcup_{n=1}^{\infty} \left(\frac{1}{n+1}, \frac{1}{n}\right]\right) \cup (2, 3)$$

$$X = [0, 1] \cup (2, 3).$$

Intuition tells us that only in the third of these decompositions have we really broken up  $X$  into its connected pieces. What distinguishes the third from the other two is that neither of the pieces  $[0, 1]$  or

Table 7. Separation properties of some useful spaces.

Space	$T_0$	$T_1$	$T_2$
Discrete	✓	✓	✓
Indiscrete	×	×	×
$\mathbb{R}$ , standard	✓	✓	✓
left/right ray	✓	×	×
Infinite cofinite	✓	✓	×
Uncountable cocountable	✓	✓	×
$x$ -inclusion/exclusion	✓	×	×
$A$ -inclusion/exclusion	×	×	×

(2, 3) can be enlarged into bigger connected subsets of  $X$ .

As connected spaces, the empty set and the singleton are considered to be *degenerate* and any connected subspace with more than one point is *non-degenerate*. At the opposite extreme of the largest possible component of a space  $X$  which is  $X$  itself, are the singletons  $\{x\}$  for every  $x \in X$ . This leads to the extremely important notion of a

**Definition A.3.3** (Totally disconnected space). A space  $X$  is totally disconnected if every pair of distinct points in it can be separated by a disconnection of  $X$ .

$X$  is totally disconnected iff the components in  $X$  are single points with the only nonempty connected subsets of  $X$  being the one-point sets: If  $x \neq y \in A \subseteq X$  are distinct points of a subset  $A$  of  $X$  then  $A = (A \cap H_1) \cup (A \cap H_2)$ , where  $X = H_1 \cup H_2$  with  $x \in H_1$  and  $y \in H_2$  is a disconnection of  $X$  (it is possible to choose  $H_1$  and  $H_2$  in this manner because  $X$  is assumed to be totally disconnected), is a separation of  $A$  that demonstrates that any subspace of a totally disconnected space with more than one point is disconnected.

A totally disconnected space has interesting physically appealing separation properties in terms of the (separated) Hausdorff spaces; here a topological space  $X$  is *Hausdorff*, or  $T_2$ , iff each two distinct points of  $X$  can be *separated* by disjoint neighborhoods, so that for every  $x \neq y \in X$ , there are neighborhoods  $M \in \mathcal{N}_x$  and  $N \in \mathcal{N}_y$  such that  $M \cap N = \emptyset$ . This means that for any two distinct points  $x \neq y \in X$ , it is impossible to find points that are arbitrarily close to both of them. Among the

properties of Hausdorff spaces, the following need to be mentioned.

- (H1)  $X$  is Hausdorff iff for each  $x \in X$  and any point  $y \neq x$ , there is a neighborhood  $N$  of  $x$  such that  $y \notin \text{Cl}(N)$ . This leads to the significant result that for any  $x \in X$  the closed singleton

$$\{x\} = \bigcap_{N \in \mathcal{N}_x} \text{Cl}(N)$$

is the intersection of the closures of any local base at that point, which in the language of nets and filters (Appendix A.1) means that a net in a Hausdorff space cannot converge to more than one point in the space and the adherent set  $\text{adh}(\mathcal{N}_x)$  of the neighborhood filter at  $x$  is the singleton  $\{x\}$ .

- (H2) Since each singleton is a closed set, each finite set in a Hausdorff space is also closed in  $X$ . Unlike a cofinite space, however, there can clearly be infinite closed sets in a Hausdorff space.
- (H3) Any point  $x$  in a Hausdorff space  $X$  is a cluster point of  $A \subseteq X$  iff every neighborhood of  $x$  contains infinitely many points of  $A$ , a fact that has led to our mental conditioning of the points of a (Cauchy) sequence piling up in neighborhoods of the limit. Thus suppose for the sake of argument that although some neighborhood of  $x$  contains only a finite number of points,  $x$  is nonetheless a cluster point of  $A$ . Then there is an open neighborhood  $U$  of  $x$  such that  $U \cap (A - \{x\}) = \{x_1, \dots, x_n\}$  is a finite closed set of  $X$  not containing  $x$ , and  $U \cap (X - \{x_1, \dots, x_n\})$  being the intersection of two open sets, is an open neighborhood of  $x$  not intersecting  $A - \{x\}$  implying thereby that  $x \notin \text{Der}(A)$ ; in fact  $U \cap (X - \{x_1, \dots, x_n\})$  is simply  $\{x\}$  if  $x \in A$  or belongs to  $\text{Bdy}_{X-A}(A)$  when  $x \in X - A$ . Conversely if every neighborhood of a point of  $X$  intersects  $A$  in infinitely many points, that point must belong to  $\text{Der}(A)$  by definition.

Weaker separation axioms than Hausdorffness are those of  $T_0$ , respectively  $T_1$ , spaces in which for every pair of distinct points *at least one*, respectively *each one*, has some neighborhood not containing the other; the following table is a listing of the separation properties of some useful spaces.

It should be noted that that as none of the properties (H1)–(H3) need neighborhoods of both points simultaneously, it is sufficient for  $X$  to be  $T_1$  for the conclusions to remain valid.

From its definition it follows that any totally disconnected space is a Hausdorff space and is therefore both  $T_1$  and  $T_0$  spaces as well. However, if a Hausdorff space has a base of clopen sets then it is totally disconnected; this is so because if  $x$  and  $y$  are distinct points of  $X$ , then the assumed property of  $x \in H \subseteq M$  for every  $M \in \mathcal{N}_x$  and some clopen set  $M$  yields  $X = H \cup (X - H)$  as a disconnection of  $X$  that separates  $x$  and  $y \in X - H$ ; note that the assumed Hausdorffness of  $X$  allows  $M$  to be chosen so as not to contain  $y$ .

**Example A.3.1**

- (1) Every indiscrete space is connected; every subset of an indiscrete space is connected. Hence if  $X$  is empty or a singleton, it is connected. A discrete space is connected iff it is either empty or is a singleton; the only connected subsets in a discrete space are the degenerate ones. This is an extreme case of lack of connectedness, and a discrete space is the simplest example of a total disconnected space.
- (2)  $\mathbb{Q}$ , the set of rationals considered as a subspace of the real line, is (totally) disconnected because all rationals larger than a given irrational  $r$  is a clopen set in  $\mathbb{Q}$ , and

$$\mathbb{Q} = \left( (-\infty, r) \cap \mathbb{Q} \right) \cup \left( \mathbb{Q} \cap (r, \infty) \right)$$

$r$  is an irrational

is the union of two disjoint clopen sets in the relative topology of  $\mathbb{Q}$ . The sets  $(-\infty, r) \cap \mathbb{Q}$  and  $\mathbb{Q} \cap (r, \infty)$  are clopen in  $\mathbb{Q}$  because neither contains a cluster point of the other. Thus for example, any neighborhood of the second must contain the irrational  $r$  in order to be able to cut the first which means that any neighborhood of a point in either of the relatively open sets cannot be wholly contained in the other. The only connected sets of  $\mathbb{Q}$  are one point subsets consisting of the individual rationals. In fact, a connected piece of  $\mathbb{Q}$ , being a connected subset of  $\mathbb{R}$ , is an interval in  $\mathbb{R}$ , and a nonempty interval cannot be contained in  $\mathbb{Q}$  unless it is a singleton. It needs to be noted that the individual points of the rational line are not (cl)open because any open subset of  $\mathbb{R}$  that contains a rational must also contain others different from

it. This example shows that a space need not be discrete for each of its points to be a component and thereby for the space to be totally disconnected.

In a similar fashion, the set of irrationals is (totally) disconnected because all the irrationals larger than a given rational is an example of a clopen set in  $\mathbb{R} - \mathbb{Q}$ .

- (3) The  $p$ -inclusion ( $A$ -inclusion) topology is connected; a subset in this topology is connected iff it is degenerate or contains  $p$ . For, a subset inherits the discrete topology if it does not contain  $p$ , and  $p$ -inclusion topology if it contains  $p$ .
- (4) The cofinite (cocountable) topology on an infinite (uncountable) space is connected; a subset in a cofinite (cocountable) space is connected iff it is degenerate or infinite (countable).
- (5) Removal of a single point may render a connected space disconnected and even totally disconnected. In the former case, the point removed is called a *cut point* and in the second, it is a *dispersion point*. Any real number is a cut point of  $\mathbb{R}$  and it does not have any dispersion point only.
- (6) Let  $X$  be a topological space. Considering components of  $X$  as equivalence classes by the equivalence relation  $\sim$  with  $Q : X \rightarrow X/\sim$  denoting the quotient map,  $X/\sim$  is totally disconnected: As  $Q^{-}([x])$  is connected for each  $[x] \in X/\sim$  in a component class of  $X$ , and as any open or closed subset  $A \subseteq X/\sim$  is connected iff  $Q^{-}(A)$  is open or closed, it must follow that  $A$  can only be a singleton.

The next notion of compactness in topological spaces provides an insight of the role of non-empty adherent sets of filters that lead in a natural fashion to the concept of attractors in the dynamical systems theory that we take up next.

**Definition A.3.4** (Compactness). A topological space  $X$  is compact iff every open cover of  $X$  contains a finite subcover of  $X$ .

This definition of compactness has an useful equivalent contrapositive reformulation: For any given collection of open sets of  $X$  if none of its finite subcollections cover  $X$ , then the entire collection also cannot cover  $X$ . The following theorem is a statement of the fundamental property of compact spaces in terms of adherences of filters in such

spaces, the proof of which uses this contrapositive characterization of compactness.

**Theorem A.2.1.** A topological space  $X$  is compact iff each class of closed subsets of  $X$  with finite intersection property has non-empty intersection.

*Proof. Necessity.* Let  $X$  be a compact space. Let  $\mathcal{F} = \{F_\alpha\}_{\alpha \in \mathbb{D}}$  be a collection of closed subsets of  $X$  with finite FIP, and let  $\mathcal{G} = \{X - F_\alpha\}_{\alpha \in \mathbb{D}}$  be the corresponding open sets of  $X$ . If  $\{G_i\}_{i=1}^N$  is a non-empty finite subcollection from  $\mathcal{G}$ , then  $\{X - G_i\}_{i=1}^N$  is the corresponding non-empty finite subcollection of  $\mathcal{F}$ . Hence from the assumed finite intersection property of  $\mathcal{F}$ , it must be true that

$$X - \bigcup_{i=1}^N G_i = \bigcap_{i=1}^N (X - G_i) \quad (\text{DeMorgan's Law}) \\ \neq \emptyset,$$

so that no finite subcollection of  $\mathcal{G}$  can cover  $X$ . Compactness of  $X$  now implies that  $\mathcal{G}$  too cannot cover  $X$  and therefore

$$\bigcap_{\alpha} F_\alpha = \bigcap_{\alpha} (X - G_\alpha) = X - \bigcup_{\alpha} G_\alpha \neq \emptyset.$$

The proof of the converse is a simple exercise of reversing the arguments involving the two equations in the proof above. ■

Our interest in this theorem and its proof lies in the following corollary — *which essentially means that for every filter  $\mathcal{F}$  on a compact space the adherent set  $\text{adh}(\mathcal{F})$  is not empty* — from which it follows that every net in a compact space must have a convergent subnet.

**Corollary.** A space  $X$  is compact iff for every class  $\mathcal{A} = (A_\alpha)$  of nonempty subsets of  $X$  with FIP,  $\text{adh}(\mathcal{A}) = \bigcap_{A_\alpha \in \mathcal{A}} \text{Cl}(A_\alpha) \neq \emptyset$ .

The proof of this result for nets given by the next theorem illustrates the general approach in such cases which is all that is basically needed in dealing with attractors of dynamical systems; compare Theorem A.1.3.

**Theorem A.3.2.** A topological space  $X$  is compact iff each net in  $X$  adheres to  $X$ .

*Proof. Necessity.* Let  $X$  be a compact space,  $\chi : \mathbb{D} \rightarrow X$  a net in  $X$ , and  $\mathbb{R}_\alpha$  the residual of  $\alpha$  in the directed set  $\mathbb{D}$ . For the filter-base  $({}_F\mathcal{B}_{\chi}(\mathbb{R}_\alpha))_{\alpha \in \mathbb{D}}$  of nonempty, decreasing, nested subsets of  $X$  associated with the net  $\chi$ , compactness of  $X$  requires from



$\bigcap_{\alpha \preceq \delta} \text{Cl}(\chi(\mathbb{R}_\alpha) \supseteq \chi(\mathbb{R}_\delta) \neq \emptyset$ , that the uncountably intersecting subset

$$\text{adh}({}_F\mathcal{B}_\chi) := \bigcap_{\alpha \in \mathbb{D}} \text{Cl}(\chi(\mathbb{R}_\alpha))$$

of  $X$  be non-empty. If  $x \in \text{adh}({}_F\mathcal{B}_\chi)$  then because  $x$  is in the closure of  $\chi(\mathbb{R}_\beta)$ , it follows from Eq. (20) that  $N \cap \chi(\mathbb{R}_\beta) \neq \emptyset$ <sup>33</sup> for every  $N \in \mathcal{N}_x$ ,  $\beta \in \mathbb{D}$ . Hence  $\chi(\gamma) \in N$  for some  $\gamma \succeq \beta$  so that  $x \in \text{adh}(\chi)$ ; see Eq. (A.16).

*Sufficiency.* Let  $\chi$  be a net in  $X$  that adheres at  $x \in X$ . From any class  $\mathcal{F}$  of closed subsets of  $X$  with FIP, construct as in the proof of Theorem A.1.4, a decreasing nested sequence of closed subsets  $C_\beta = \bigcap_{\alpha \preceq \beta \in \mathbb{D}} \{F_\alpha : F_\alpha \in \mathcal{F}\}$  and consider the directed set  ${}_{\mathbb{D}}C_\beta = \{(C_\beta, \beta) : (\beta \in \mathbb{D})(x_\beta \in C_\beta)\}$  with its natural direction (A.11) to define the net  $\chi(C_\beta, \beta) = x_\beta$  in  $X$ ; see Definition A.1.10. From the assumed adherence of  $\chi$  at some  $x \in X$ , it follows that  $N \cap F \neq \emptyset$  for every  $N \in \mathcal{N}_x$  and  $F \in \mathcal{F}$ . Hence  $x$  belongs to the closed set  $F$  so that  $x \in \text{adh}(\mathcal{F})$ ; see Eq. (A.24). Hence  $X$  is compact. ■

Using Theorem A.1.5 that specifies a definite criterion for the adherence of a net, this theorem reduces to the useful formulation that a space is compact iff each net in it has some convergent subnet. An important application is the following: Since every decreasing sequence  $(F_m)$  of nonempty sets has FIP (because  $\bigcap_{m=1}^M F_m = F_M$  for every finite  $M$ ), every decreasing sequence of nonempty closed subsets of a compact space has nonempty intersection. For a complete metric space this is known as the *Nested Set Theorem*, and for  $[0, 1]$  and other compact subspaces of  $\mathbb{R}$  as the *Cantor Intersection Theorem*.<sup>34</sup>

For subspaces  $A$  of  $X$ , it is the relative topology that determines as usual compactness of  $A$ ; however the following criterion renders this test in terms of the relative topology unnecessary and shows that the topology of  $X$  itself is sufficient to determine

compactness of subspaces: *A subspace  $K$  of a topological space  $X$  is compact iff each open cover of  $K$  in  $X$  contains a finite cover of  $K$ .*

A proper understanding of the distinction between compactness and closedness of subspaces — which often causes much confusion to the non-specialist — is expressed in the next two theorems. As a motivation for the first that establishes that not every subset of a compact space need be compact, mention may be made of the subset  $(a, b)$  of the compact closed interval  $[a, b]$  in  $\mathbb{R}$ .

**Theorem A.3.3.** *A closed subset  $F$  of a compact space  $X$  is compact.*

*Proof.* Let  $\mathcal{G}$  be an open cover of  $F$  so that an open cover of  $X$  is  $\mathcal{G} \cup (X - F)$ , which because of compactness of  $X$  contains a finite subcover  $\mathcal{U}$ . Then  $\mathcal{U} - (X - F)$  is a finite collection of  $\mathcal{G}$  that covers  $F$ . ■

It is not true in general that a compact subset of a space is necessarily closed. For example, in an infinite set  $X$  with the cofinite topology, let  $F$  be an infinite subset of  $X$  with  $X - F$  also infinite. Then although  $F$  is not closed in  $X$ , it is nevertheless compact because  $X$  is compact. Indeed, let  $\mathcal{G}$  be an open cover of  $X$  and choose any non-empty  $G_0 \in \mathcal{G}$ . If  $G_0 = X$  then  $\{G_0\}$  is the required finite cover of  $X$ . If this is not the case, then because  $X - G_0 = \{x_i\}_{i=1}^n$  is a finite set, there is a  $G_i \in \mathcal{G}$  with  $x_i \in G_i$  for each  $1 \leq i \leq n$ , and therefore  $\{G_i\}_{i=0}^n$  is the finite cover that demonstrates the compactness of the cofinite space  $X$ . Compactness of  $F$  now follows because the subspace topology on  $F$  is the induced cofinite topology from  $X$ . The distinguishing feature of this topology is that it, like the cocountable, is not Hausdorff: If  $U$  and  $V$  are any two nonempty open sets of  $X$ , then they cannot be disjoint as the complements of the open sets can only be finite and if  $U \cap V$  were to be indeed empty, then

$$X = X - \emptyset = X - (U \cap V) = (X - U) \cup (X - V)$$

<sup>33</sup>This is of course a triviality if we identify each  $\chi(\mathbb{R}_\beta)$  (or  $F$  in the proof of the converse that follows) with a neighborhood  $N$  of  $X$  that generates a topology on  $X$ .

<sup>34</sup>**Nested-set theorem.** *If  $(E_n)$  is a decreasing sequence of nonempty, closed, subsets of a complete metric space  $(X, d)$  such that  $\lim_{n \rightarrow \infty} \text{dia}(E_n) = 0$ , then there is a unique point*

$$x \in \bigcap_{n=0}^{\infty} E_n.$$

The uniqueness arises because the limiting condition on the diameters of  $E_n$  imply, from property (H1), that  $(X, d)$  is a Hausdorff space.

would be a finite set. An immediate fallout of this is that in an infinite cofinite space, a sequence  $(x_i)_{i \in \mathbb{N}}$  (and even a net) with  $x_i \neq x_j$  for  $i \neq j$  behaves in an extremely unusual way: *It converges, as in the indiscrete space, to every point of the space.* Indeed if  $x \in X$ , where  $X$  is an infinite set provided with its cofinite topology, and  $U$  is any neighborhood of  $x$ , any infinite sequence  $(x_i)_{i \in \mathbb{N}}$  in  $X$  must be eventually in  $U$  because  $X - U$  is finite, and ignoring of the initial set of its values lying in  $X - U$  in no way alters the ultimate behavior of the sequence (note that this implies that the filter induced on  $X$  by the sequence agrees with its topology). Thus  $x_i \rightarrow x$  for any  $x \in X$  is a reflection of the fact that there are no small neighborhoods of any point of  $X$  with every neighborhood being almost the whole of  $X$ , except for a null set consisting of only a finite number of points. This is in sharp contrast with Hausdorff spaces where, although every finite set is also closed, every point has arbitrarily small neighborhoods that lead to unique limits of sequences. A corresponding result for cocountable spaces can be found in Example A.1.2, continued.

This example of the cofinite topology motivates the following “converse” of the previous theorem.

**Theorem A.3.4.** *Every compact subspace of a Hausdorff space is closed.*

*Proof.* Let  $K$  be a non-empty compact subset of  $X$ , Fig. 23, and let  $x \in X - K$ . Because of the separation of  $X$ , for every  $y \in K$  there are disjoint open subsets  $U_y$  and  $V_y$  of  $X$  with  $y \in U_y$ , and  $x \in V_y$ . Hence  $\{U_y\}_{y \in K}$  is an open cover for  $K$ , and from its compactness there is a finite subset  $A$  of  $K$  such that  $K \subseteq \bigcup_{y \in A} U_y$  with  $V = \bigcap_{y \in A} V_y$  an open neighborhood of  $x$ ;  $V$  is open because each

$V_y$  is a neighborhood of  $x$  and the intersection is over finitely many points  $y$  of  $A$ . To prove that  $K$  is closed in  $X$  it is enough to show that  $V$  is disjoint from  $K$ : If there is indeed some  $z \in V \cap K$  then  $z$  must be in some  $U_y$  for  $y \in A$ . But as  $z \in V$  it is also in  $V_y$  which is impossible as  $U_y$  and  $V_y$  are to be disjoint. This last part of the argument in fact shows that *if  $K$  is a compact subspace of a Hausdorff space  $X$  and  $x \notin K$ , then there are disjoint open sets  $U$  and  $V$  of  $X$  containing  $x$  and  $K$ .* ■

The last two theorems may be combined to give the obviously important

**Corollary.** *In a compact Hausdorff space, closedness and compactness of its subsets are equivalent concepts.*

In the absence of Hausdorffness, it is not possible to conclude from the assumed compactness of the space that every point to which the net may converge actually belongs to the subspace.

**Definition A.3.5.** A subset  $D$  of a topological space  $(X, \mathcal{U})$  is dense in  $X$  if  $\text{Cl}(D) = X$ . Thus the closure of  $D$  is the largest open subset of  $X$ , and every neighborhood of any point of  $X$  contains a point of  $D$  not necessarily distinct from it; refer to the distinction between Eqs. (20) and (22).

Loosely,  $D$  is dense in  $X$  iff every point of  $X$  has points of  $D$  arbitrarily close to it. A *self-dense* (dense in itself) set is a set without any isolated points; hence  $A$  is self-dense iff  $A \subseteq \text{Der}(A)$ . A closed self-dense set is called a *perfect set* so that a closed set  $A$  is perfect iff it has no isolated points. Accordingly

$$A \text{ is perfect} \Leftrightarrow A = \text{Der}(A),$$

means that the closure of a set without any isolated points is a perfect set.

**Theorem A.3.5.** *The following are equivalent statements.*

- (1)  $D$  is dense in  $X$ .
- (2) If  $F$  is any closed set of  $X$  with  $D \subseteq F$ , then  $F = X$ ; thus the only closed superset of  $D$  is  $X$ .
- (3) Every nonempty (basic) open set of  $X$  cuts  $D$ ; thus the only open set disjoint from  $D$  is the empty set  $\emptyset$ .
- (4) The exterior of  $D$  is empty.

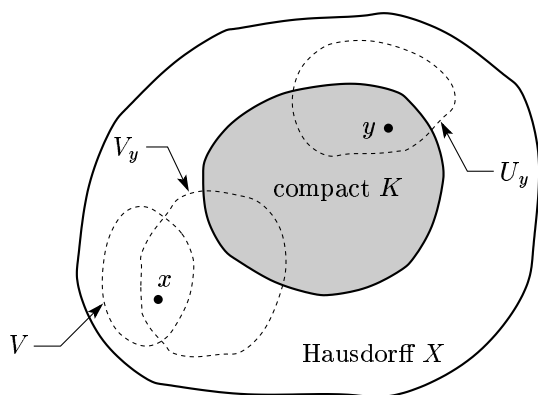


Fig. 23. Closedness of compact subsets of a Hausdorff space.

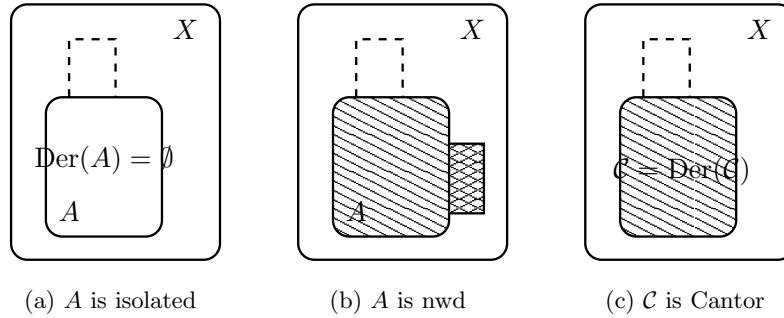


Fig. 24. Shows the distinction between isolated, nowhere dense and Cantor sets. Topologically, the Cantor set can be described as a perfect, nowhere dense, totally disconnected and compact subset of a space. (b) The closed nowhere dense set  $\text{Cl}(A)$  is the boundary of its open complement. Here downward and upward inclined hatching denote respectively  $\text{Bdy}_A(X - A)$  and  $\text{Bdy}_{X-A}(A)$ .

*Proof.* (3) If  $U$  indeed is a non-empty open set of  $X$  with  $U \cap D = \emptyset$ , then  $D \subseteq X - U \neq X$  leads to the contradiction  $X = \text{Cl}(D) \subseteq \text{Cl}(X - U) = X - U \neq X$ , which also incidentally proves (2). From (3) it follows that for any open set  $U$  of  $X$ ,  $\text{Cl}(U) = \text{Cl}(U \cap D)$  because if  $V$  is any open neighborhood of  $x \in \text{Cl}(U)$  then  $V \cap U$  is a non-empty open set of  $X$  that must cut  $D$  so that  $V \cap (U \cap D) \neq \emptyset$  implies  $x \in \text{Cl}(U \cap D)$ . Finally,  $\text{Cl}(U \cap D) \subseteq \text{Cl}(U)$  completes the proof. ■

**Definition A.3.6.** (a) A set  $A \subseteq X$  is said to be nowhere dense in  $X$  if  $\text{Int}(\text{Cl}(A)) = \emptyset$  and residual in  $X$  if  $\text{Int}(A) = \emptyset$ .

$A$  is nowhere dense in  $X$  iff

$$\text{Bdy}(X - \text{Cl}(A)) = \text{Bdy}(\text{Cl}(A)) = \text{Cl}(A)$$

so that

$$\text{Cl}(X - \text{Cl}(A)) = (X - \text{Cl}(A)) \cup \text{Cl}(A) = X$$

from which it follows that

$$A \text{ is nwd in } X \Leftrightarrow X - \text{Cl}(A) \text{ is dense in } X$$

and

$$A \text{ is residual in } X \Leftrightarrow X - A \text{ is dense in } X.$$

Thus  $A$  is nowhere dense iff  $\text{Ext}(A) := X - \text{Cl}(A)$  is dense in  $X$ , and in particular, a closed set is nowhere dense in  $X$  iff its complement is open dense in  $X$  with open-denseness being complementarily dual to closed-nowhere denseness. The rationals in reals is an example of a set that is residual but not nowhere dense. The following are readily verifiable properties of subsets of  $X$ .

- (1) A set  $A \subseteq X$  is nowhere dense in  $X$  iff it is contained in its own boundary, iff it is contained in the closure of the complement of its

closure, that is  $A \subseteq \text{Cl}(X - \text{Cl}(A))$ . In particular a closed subset  $A$  is nowhere dense in  $X$  iff  $A = \text{Bdy}(A)$ , that is iff it contains no open set.

- (2) From  $M \subseteq N \Rightarrow \text{Cl}(M) \subseteq \text{Cl}(N)$  it follows, with  $M = X - \text{Cl}(A)$  and  $N = X - A$ , that a nowhere dense set is residual, but a residual set need not be nowhere dense unless it is also closed in  $X$ .
- (3) Since  $\text{Cl}(\text{Cl}(A)) = \text{Cl}(A)$ ,  $\text{Cl}(A)$  is nowhere dense in  $X$  iff  $A$  is.
- (4) For any  $A \subseteq X$ , both  $\text{Bdy}_A(X - A) := \text{Cl}(X - A) \cap A$  and  $\text{Bdy}_{X-A}(A) := \text{Cl}(A) \cap (X - A)$  are residual sets and as Fig. 22 shows  $\text{Bdy}_X(A) = \text{Bdy}_{X-A}(A) \cup \text{Bdy}_A(X - A)$  is the union of these two residual sets. When  $A$  is closed (or open) with  $X$  its boundary, consisting of the only component  $\text{Bdy}_A(X - A)$  (or  $\text{Bdy}_{X-A}(A)$ ) as shown by the second row (or column) of the figure, being a closed set of  $X$  is also nowhere dense in  $X$ ; in fact a closed nowhere dense set is always the boundary of some open set. Otherwise, the boundary components of the two residual parts — as in the donor-donor, donor-neutral, neutral-donor and neutral-neutral cases — need not be individually closed in  $X$  (although their union is) and their union is a residual set that need not be nowhere dense in  $X$ : the union of two nowhere dense sets is nowhere dense but the union of a residual and a nowhere dense set is a residual set. One way in which a two-component boundary can be nowhere dense is by having  $\text{Bdy}_A(X - A) \supseteq \text{Der}(A)$  or  $\text{Bdy}_{X-A}(A) \supseteq \text{Der}(X - A)$ , so that it is effectively in one piece rather than in two, as show in Fig. 24(b).

**Theorem A.3.6.** *A is nowhere dense in X iff each non-empty open set of X has a non-empty open subset disjoint from Cl(A).*

*Proof.* If  $U$  is a non-empty open set of  $X$ , then  $U_0 = U \cap \text{Ext}(A) \neq \emptyset$  as  $\text{Ext}(A)$  is dense in  $X$ ;  $U_0$  is the open subset that is disjoint from  $\text{Cl}(A)$ . It clearly follows from this that each non-empty open set of  $X$  has a non-empty open subset disjoint from a nowhere dense set  $A$ . ■

What this result (which follows just from the definition of nowhere dense sets) actually means is that no point in  $\text{Bdy}_{X-A}(A)$  can be isolated in it.

**Corollary.** *A is nowhere dense in X iff Cl(A) does not contain any non-empty open set of X iff any nonempty open set that contains A also contains its closure.*

**Example A.3.2.** Each finite subset of  $\mathbb{R}^n$  is nowhere dense in  $\mathbb{R}^n$ ; the set  $\{1/n\}_{n=1}^\infty$  is nowhere dense in  $\mathbb{R}$ . The Cantor set  $\mathcal{C}$  is nowhere dense in  $[0, 1]$  because every neighborhood of any point in  $\mathcal{C}$  must contain, by its very construction, a point with 1 in its ternary representation. That the interior and the interior of the closure of a set are not necessarily the same is seen in the example of the rationals in reals: The set of rational numbers  $\mathbb{Q}$  has empty interior because any neighborhood of a rational number contains irrational numbers (so also is the case for irrational numbers) and  $\mathbb{R} = \text{Int}(\text{Cl}(\mathbb{Q})) \supseteq \text{Int}(\mathbb{Q}) = \emptyset$  justifies the notion of a nowhere dense set.

The following properties of  $\mathcal{C}$  can be taken to define any subset of a topological space as a Cantor set; set-theoretically it should be clear from its classical middle-third construction that the Cantor set consists of all points of the closed interval  $[0, 1]$  whose infinite triadic (base 3) representation, expressed so as not to terminate with an infinite string of 1's, does not contain the digit 1. Accordingly, any end point of the infinite set of closed intervals whose intersection yields the Cantor set, is represented by a repeating string of either 0 or 2 while a non end point has every other arbitrary collection of these two digits. Recalling that any number in  $[0, 1]$  is a rational iff its representation in any base is terminating or recurring — thus any decimal that neither repeats or terminates but consists of all possible sequences of all possible digits represents an irrational

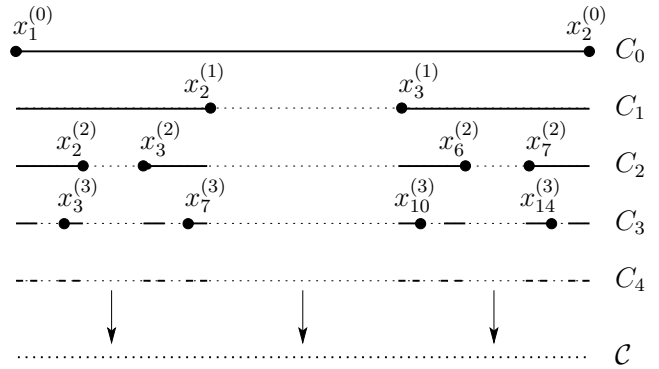


Fig. 25. Construction of the classical 1/3-Cantor set. The end points of  $C_3$ , for example, in increasing order are:  $|0, \frac{1}{27}|; |\frac{2}{27}, \frac{1}{9}|; |\frac{2}{9}, \frac{7}{27}|; |\frac{8}{27}, \frac{1}{3}|; |\frac{2}{3}, \frac{19}{27}|; |\frac{20}{27}, \frac{7}{9}|; |\frac{8}{9}, \frac{25}{27}|; |\frac{26}{27}, 1|$ .  $C_i$  is the union of  $2^i$  pairwise disjoint closed intervals each of length  $3^{-i}$  and the non-empty infinite intersection  $\mathcal{C} = \cap_{i=0}^\infty C_i$  is the adherent Cantor set of the filter-base of closed sets  $\{C_0, C_1, C_2, \dots\}$ .

number — it follows that both rationals and irrationals belong to the Cantor set.

(C1)  **$\mathcal{C}$  is totally disconnected.** If possible, let  $\mathcal{C}$  have a component containing points  $a$  and  $b$  with  $a < b$ . Then  $[a, b] \subseteq \mathcal{C} \Rightarrow [a, b] \subseteq C_i$  for all  $i$ . But this is impossible because we may choose  $i$  large enough to have  $3^{-i} < b - a$  so that  $a$  and  $b$  must belong to two different members of the pairwise disjoint closed  $2^i$  subintervals each of length  $3^{-i}$  that constitutes  $C_i$ . Hence

$$[a, b] \text{ is not a subset of any } C_i \Rightarrow [a, b] \text{ is not a subset of } \mathcal{C}.$$

(C2)  **$\mathcal{C}$  is perfect** so that for any  $x \in \mathcal{C}$  every neighborhood of  $x$  must contain some other point of  $\mathcal{C}$ . Supposing to the contrary that the singleton  $\{x\}$  is an open set of  $\mathcal{C}$ , there must be an  $\varepsilon > 0$  such that in the usual topology of  $\mathbb{R}$

$$\{x\} = \mathcal{C} \cap (x - \varepsilon, x + \varepsilon). \tag{A.52}$$

Choose a positive integer  $i$  large enough to satisfy  $3^{-i} < \varepsilon$ . Since  $x$  is in every  $C_i$ , it must be in one of the  $2^i$  pairwise disjoint closed intervals  $[a, b] \subset (x - \varepsilon, x + \varepsilon)$  each of length  $3^{-i}$  whose union is  $C_i$ . As  $[a, b]$  is an interval, at least one of the end points of  $[a, b]$  is different from  $x$ , and since an end point belongs to  $\mathcal{C}$ ,  $\mathcal{C} \cap (x - \varepsilon, x + \varepsilon)$  must also contain this point thereby violating Eq. (A.52).

- (C3)  $\mathcal{C}$  is **nowhere dense** because each neighborhood of any point of  $\mathcal{C}$  intersects  $\text{Ext}(\mathcal{C})$ ; see Theorem A.3.6.
- (C4)  $\mathcal{C}$  is **compact** because it is a closed subset contained in the compact subspace  $[0, 1]$  of  $\mathbb{R}$ , see Theorem A.3.3. The compactness of  $[0, 1]$  follows from the Heine-Borel Theorem which states that any subset of the real line is compact iff it is both closed and bounded with respect to the Euclidean metric on  $\mathbb{R}$ .

Compare (C1) and (C2) with the essentially similar arguments of Example A.3.1(2) for the subspace of rationals in  $\mathbb{R}$ .

### A.4. Neutron Transport Theory

This section introduces the reader to the basics of the *linear* neutron transport theory where graphical convergence approximations to the singular distributions, interpreted here as multifunctions, led to the study of this paper. The one-speed (that is mono-energetic) neutron transport equation in one dimension and plane geometry, is

$$\mu \frac{\partial \Phi(x, \mu)}{\partial x} + \Phi(x, \mu) = \frac{c}{2} \int_{-1}^1 \Phi(x, \mu') d\mu', \quad 0 < c < 1, \quad -1 \leq \mu \leq 1 \tag{A.53}$$

where  $x$  is a non-dimensional physical space variable that denotes the location of the neutron moving in a direction  $\theta = \cos^{-1}(\mu)$ ,  $\Phi(x, \mu)$  is a neutron density distribution function such that  $\Phi(x, \mu) dx d\mu$  is the expected number of neutrons in a distance  $dx$  about the point  $x$  moving at constant speed with their direction cosines of motion in  $d\mu$  about  $\mu$ , and  $c$  is a physical constant that will be taken to satisfy the restriction shown above. Case's method starts by assuming the solution to be of the form  $\Phi_\nu(x, \mu) = e^{-x/\mu} \phi(\mu, \nu)$  with a normalization integral constraint of  $\int_{-1}^1 \phi(\mu, \nu) d\mu = 1$  to lead to the simple equation

$$(\nu - \mu)\phi(\mu, \nu) = \frac{c\nu}{2} \tag{A.54}$$

for the unknown function  $\phi(\nu, \mu)$ . Case then suggested, see [Case & Zweifel, 1967], the non-simple complete solution of this equation to be

$$\phi(\mu, \nu) = \frac{c\nu}{2} \mathcal{P} \frac{1}{\nu - \mu} + \lambda(\nu)\delta(\nu - \mu), \tag{A.55}$$

where  $\lambda(\nu)$  is the usual combination coefficient of the solutions of the homogeneous and non-homogeneous parts of a linear equation,  $\mathcal{P}(\cdot)$  is a principal value and  $\delta(x)$  the Dirac delta, to lead to the full-range  $-1 \leq \mu \leq 1$  solution valid for  $-\infty < x < \infty$

$$\begin{aligned} \Phi(x, \mu) &= a(\nu_0)e^{-x/\nu_0}\phi(\mu, \nu_0) \\ &+ a(-\nu_0)e^{x/\nu_0}\phi(-\nu_0, \mu) \\ &+ \int_{-1}^1 a(\nu)e^{-x/\nu}\phi(\mu, \nu)d\nu \end{aligned} \tag{A.56}$$

of the one-speed neutron transport equation (A.53). Here the real  $\nu_0$  and  $\nu$  satisfy respectively the integral constraints

$$\begin{aligned} \frac{c\nu_0}{2} \ln \frac{\nu_0 + 1}{\nu_0 - 1} &= 1, \quad |\nu_0| > 1 \\ \lambda(\nu) &= 1 - \frac{c\nu}{2} \ln \frac{1 + \nu}{1 - \nu}, \quad \nu \in [-1, 1], \end{aligned}$$

with

$$\phi(\mu, \nu_0) = \frac{c\nu_0}{2} \frac{1}{\nu_0 - \mu}$$

following from Eq. (A.55).

It can be shown [Case & Zweifel, 1967] that the eigenfunctions  $\phi(\nu, \mu)$  satisfy the full-range orthogonality condition

$$\int_{-1}^1 \mu \phi(\nu, \mu) \phi(\nu', \mu) d\mu = N(\nu) \delta(\nu - \nu'),$$

where the odd normalization constants  $N$  are given by

$$\begin{aligned} N(\pm\nu_0) &= \int_{-1}^1 \mu \phi^2(\pm\nu_0, \mu) d\mu \quad \text{for } |\nu_0| > 1 \\ &= \pm \frac{c\nu_0^3}{2} \left( \frac{c}{\nu_0^2 - 1} - \frac{1}{\nu_0^2} \right), \end{aligned}$$

and

$$N(\nu) = \nu \left( \lambda^2(\nu) + \left( \frac{\pi c \nu}{2} \right)^2 \right) \quad \text{for } \nu \in [-1, 1].$$

With a source of particles  $\psi(x_0, \mu)$  located at  $x = x_0$  in an infinite medium, Eq. (A.56) reduces to the boundary condition, with  $\mu, \nu \in [-1, 1]$ ,

$$\begin{aligned} \psi(x_0, \mu) &= a(\nu_0)e^{-x_0/\nu_0}\phi(\mu, \nu_0) \\ &+ a(-\nu_0)e^{x_0/\nu_0}\phi(-\nu_0, \mu) \\ &+ \int_{-1}^1 a(\nu)e^{-x_0/\nu}\phi(\mu, \nu)d\nu \end{aligned} \tag{A.57}$$

for the determination of the expansion coefficients  $a(\pm\nu_0)$ ,  $\{a(\nu)\}_{\nu \in [-1,1]}$ . Use of the above orthogonality integrals then lead to the complete solution of the problem to be

$$a(\nu) = \frac{e^{x_0/\nu}}{N(\nu)} \int_{-1}^1 \mu \psi(x_0, \mu) \phi(\mu, \nu) d\mu, \\ \nu = \pm\nu_0 \text{ or } \nu \in [-1, 1].$$

For example, in the infinite-medium Greens function problem with  $x_0 = 0$  and  $\psi(x_0, \mu) = \delta(\mu - \mu_0)/\mu$ , the coefficients are  $a(\pm\nu_0) = \phi(\mu_0, \pm\nu_0)/N(\pm\nu_0)$  when  $\nu = \pm\nu_0$ , and  $a(\nu) = \phi(\mu_0, \nu)/N(\nu)$  for  $\nu \in [-1, 1]$ .

For a half-space  $0 \leq x < \infty$ , the obvious reduction of Eq. (A.56) to

$$\Phi(x, \mu) = a(\nu_0)e^{-x/\nu_0}\phi(\mu, \nu_0) + \int_0^1 a(\nu)e^{-x/\nu}\phi(\mu, \nu)d\nu \quad (\text{A.58})$$

with boundary condition,  $\mu, \nu \in [0, 1]$ ,

$$\psi(x_0, \mu) = a(\nu_0)e^{-x_0/\nu_0}\phi(\mu, \nu_0) + \int_0^1 a(\nu)e^{-x_0/\nu}\phi(\mu, \nu)d\nu, \quad (\text{A.59})$$

leads to an infinitely more difficult determination of the expansion coefficients due to the more involved nature of the orthogonality relations of the eigenfunctions in the half-interval  $[0, 1]$  that now reads for  $\nu, \nu' \in [0, 1]$  [Case & Zweifel, 1967]

$$\int_0^1 W(\mu)\phi(\mu, \nu')\phi(\mu, \nu)d\mu = \frac{W(\nu)N(\nu)}{\nu}\delta(\nu - \nu') \\ \int_0^1 W(\mu)\phi(\mu, \nu_0)\phi(\mu, \nu)d\mu = 0 \\ \int_0^1 W(\mu)\phi(\mu, -\nu_0)\phi(\mu, \nu)d\mu = c\nu\nu_0X(-\nu_0)\phi(\nu, -\nu_0) \quad (\text{A.60}) \\ \int_0^1 W(\mu)\phi(\mu, \pm\nu_0)\phi(\mu, \nu_0)d\mu = \mp \left(\frac{c\nu_0}{2}\right)^2 X(\pm\nu_0) \\ \int_0^1 W(\mu)\phi(\mu, \nu_0)\phi(\mu, -\nu)d\mu = \frac{c^2\nu\nu_0}{4}X(-\nu)$$

$$\int_0^1 W(\mu)\phi(\mu, \nu')\phi(\mu, -\nu)d\mu = \frac{c\nu'}{2}(\nu_0 + \nu)\phi(\nu', -\nu)X(-\nu)$$

where the half-range weight function  $W(\mu)$  is defined as

$$W(\mu) = \frac{c\mu}{2(1-c)(\nu_0 + \mu)X(-\mu)} \quad (\text{A.61})$$

in terms of the  $X$ -function

$$X(-\mu) = \exp - \left\{ \frac{c}{2} \int_0^1 \frac{\nu}{N(\nu)} \left[ 1 + \frac{c\nu^2}{1-\nu^2} \right] \ln(\nu + \mu) d\nu \right\}, \\ 0 \leq \mu \leq 1,$$

that is conveniently obtained from a numerical solution of the nonlinear integral equation

$$\Omega(-\mu) = 1 - \frac{c\mu}{2(1-c)} \int_0^1 \frac{\nu_0^2(1-c) - \nu^2}{(\nu_0^2 - \nu^2)(\mu + \nu)\Omega(-\nu)} d\nu \quad (\text{A.62})$$

to yield

$$X(-\mu) = \frac{\Omega(-\mu)}{\mu + \nu_0\sqrt{1-c}},$$

and  $X(\pm\nu_0)$  satisfy

$$X(\nu_0)X(-\nu_0) = \frac{\nu_0^2(1-c) - 1}{2(1-c)\nu_0^2(\nu_0^2 - 1)}.$$

Two other useful relations involving the  $W$ -function are given by  $\int_0^1 W(\mu)\phi(\mu, \nu_0)d\mu = c\nu_0/2$  and  $\int_0^1 W(\mu)\phi(\mu, \nu)d\mu = c\nu/2$ .

The utility of these full- and half-range orthogonality relations lie in the fact that a suitable class of functions of the type that is involved here can always be expanded in its terms, see [Case & Zweifel, 1967]. An example of this for a full-range problem has been given above; we end this introduction to the generalized — traditionally known as singular in neutron transport theory — eigenfunction method with two examples of half-range orthogonality integrals to the half-space problems A and B of Sec. 5.

**Problem A** (The Milne Problem). In this case there is no incident flux of particles from outside the medium at  $x = 0$ , but for large  $x > 0$  the neutron distribution inside the medium behaves like  $e^{x/\nu_0}\phi(-\nu_0, \mu)$ . Hence the boundary condition

(A.59) at  $x = 0$  reduces to

$$-\phi(\mu, -\nu_0) = a_A(\nu_0)\phi(\mu, \nu_0) + \int_0^1 a_A(\nu)\phi(\mu, \nu)d\nu \quad \mu \geq 0$$

Use of the fourth and third equations of Eq. (A.60) and the explicit relation Eq. (A.61) for  $W(\mu)$  gives respectively the coefficients

$$a_A(\nu_0) = \frac{X(-\nu_0)}{X(\nu_0)}$$

$$a_A(\nu) = -\frac{1}{N(\nu)}c(1-c)\nu_0^2\nu X(-\nu_0)X(-\nu). \tag{A.63}$$

The extrapolated end point  $z_0$  of Eq. (67) is related to  $a_A(\nu_0)$  of the Milne problem by  $a_A(\nu_0) = -\exp(-2z_0/\nu_0)$ .

**Problem B** (The Constant Source Problem). Here the boundary condition at  $x = 0$  is

$$1 = a_B(\nu_0)\phi(\mu, \nu_0) + \int_0^1 a_B(\nu)\phi(\mu, \nu)d\nu \quad \mu \geq 0$$

which leads, using the integral relations satisfied by  $W$ , to the expansion coefficients

$$a_B(\nu_0) = -2/c\nu_0 X(\nu_0)$$

$$a_B(\nu) = \frac{1}{N(\nu)}(1-c)\nu(\nu_0 + \nu)X(-\nu). \tag{A.64}$$

where  $X(\pm\nu_0)$  are related to Problem A as

$$X(\nu_0) = \frac{1}{\nu_0} \sqrt{\frac{\nu_0^2(1-c) - 1}{2a_A(\nu_0)(1-c)(\nu_0^2 - 1)}}$$

$$X(-\nu_0) = \frac{1}{\nu_0} \sqrt{\frac{a_A(\nu_0)(\nu_0^2(1-c) - 1)}{2(1-c)(\nu_0^2 - 1)}}.$$

This brief introduction to the singular eigenfunction method should convince the reader of the great difficulties associated with half-space, half-range methods in particle transport theory; note that the  $X$ -functions in the coefficients above must be obtained from numerically computed tables. In contrast, full-range methods are more direct due to the simplicity of the weight function  $\mu$ , which suggests the full-range formulation of half-range problems presented in Sec. 5. Finally it should be mentioned that this singular eigenfunction method is based on the theory of singular integral equations.