

Digital Geometry and Mathematical Morphology

Lecture Notes, Spring Semester 2002

Christer O. Kiselman

Contents:

- 1. Introduction**
 - 1.1. Why digital geometry?
 - 1.2. Why mathematical morphology?
- 2. Morphological operations on sets and functions**
 - 2.1. Infimal convolution
 - 2.2. Dilations and erosions
 - 2.3. Exercises
- 3. Ordered sets: first notions**
- 4. Closure operators on subsets of an abelian group**
 - 4.1. Combining erosions and dilations
 - 4.2. Characterizing closure operators which commute with translations
 - 4.3. Matheron's structure theorems
 - 4.4. Exercises
- 5. Distance transforms**
 - 5.1. Definition and basic properties of distance transforms
 - 5.2. Distance transforms and sublevel sets
 - 5.3. Distance transforms in normed vector spaces
 - 5.4. Chamfer distances
 - 5.5. Comparing distances
 - 5.6. The calculus of balls
 - 5.7. Exercises
- 6. Skeletonizing**
 - 6.1. Definition of the skeleton
 - 6.2. Existence of skeletons
 - 6.3. Properties of skeletons
 - 6.4. Exercises
- 7. Lattices**
 - 7.1. Definition and first properties of lattices
 - 7.2. Morphology on lattices
 - 7.3. Inverses of mappings between lattices
 - 7.4. Division of mappings between lattices
 - 7.5. Exercises
- 8. Notions of topology**
 - 8.1. Mappings
 - 8.2. Definition of topologies
 - 8.3. Transport of topologies
 - 8.4. Continuous mappings
 - 8.5. Connectedness
 - 8.6. Quotient spaces
 - 8.7. Separation axioms
 - 8.8. Smallest neighborhood spaces
 - 8.9. Exercises

Contents, cont'd:

9. A closer look at the Khalimsky plane

- 9.1. Continuous functions
- 9.2. A fixed-point theorem
- 9.3. Jordan curve theorems
- 9.4. Exercises

10. Digitization

- 10.1. What is digitization?
- 10.2. Voronoi cells
- 10.3. Digital lines
- 10.4. Exercises

References

1. Introduction

Digital geometry is, simply put, the geometry of the computer screen. Mathematical morphology is, in equally simple words, the theory and practice of transformations of sets and functions with an emphasis on their shape. In many cases these transformations have been known for a long time, but they have come into focus for the same reason as digital geometry: the operations can actually be performed on a computer.

These notes form the basis of my lectures at Uppsala University during the spring semester of 2002. The reader will probably notice that I am trying to learn the subject at the same time. Hopefully this will not make the notes more difficult to read.

1.1. Why digital geometry?

As everyone knows by now, there is a universal answer (*Why not?*) to all such questions, but I shall nevertheless try to give some motivation for this relatively new field and why I think it is worthwhile as a mathematical theory.

Points, straight lines and planes have been studied for well over two millenia, and certain curves, like ellipses and hyperbolas, have been an object for our curiosity for almost as long. More general curves, like lemniscates and cardioids, have been studied for several centuries. In the study of these curves we rely very much on the fact that we can draw them on paper. But with the advent of computers we have acquired a new method of drawing pictures. On the computer screen we see images, and the images consist of little picture elements, *pixels*, that the eye puts together to form geometric objects. A straight line is therefore not what Euclid understood by a straight line, but rather a finite collection of dots on the screen, which the eye nevertheless perceives as a connected line segment. Is there a geometry for these images on the screen? The answer is in the affirmative. We need not be content with the images as more or less accurate approximations of ideal straight lines or curves: we can treat these finite sets of points with the same accuracy as that which Euclid had in his geometry. This is digital geometry.

The field is young in comparison with Euclid's: the notion of a straight line was clarified in 1974 by Azriel Rosenfeld. We can also talk about curves in the digital plane. In fact, we can take any notion in Euclidean geometry, try to translate it to digital geometry, and see if a certain result in Euclidean geometry becomes true in a digital interpretation. The Jordan curve theorem is an instructive example.

For an elementary discussion of digital geometry and in particular of the Jordan curve theorem, see Kiselman (MS).

The plane may be divided into triangles, rectangles, or hexagons. Then the centers of the pixels form, respectively, a hexagonal, rectangular, or triangular pattern (Figures 1 and 2). In all these cases we may use a pair of integer coordinates to indicate the location of a pixel. This is obvious and easy in the rectangular case, but useful also in the other two cases, although we need to be careful about the metric then. Therefore we often speak of \mathbf{Z}^2 as the set of pixels, although, more precisely, a pair of integers $x = (x_1, x_2) \in \mathbf{Z}^2$ is just the address of a pixel.

1.2. Why mathematical morphology?

The field that has become known as mathematical morphology is quite old in a sense; it is about operations on sets and functions that have been around for a long time, but which are now being systematized and studied under a new angle, precisely because it is possible to actually perform operations on the computer and see on the screen what happens.

My personal view is that morphology has its origin in our trying to understand a complicated world. The world is so complex that the human mind—and the human eye—cannot perceive all its minute details, but needs a simplified image, a simplified structure. The need to simplify a complicated object is, in this view of things, the basic impulse behind mathematical morphology, and this is what mathematical morphology does. Related to this is the fact that an image may contain a lot of disturbances, or rather, it almost always does. Therefore, most images need to be tidied up. Thus another need to process images; it is related to the first, for the border line between dirt and other kind of disturbances is not too clear.

Let us think first of Euclidean geometry, and consider cardinalities. The set \mathbf{N} of nonnegative integers is infinite, and its cardinality is denoted by $\text{card}(\mathbf{N}) = \aleph_0$. The set of real numbers \mathbf{R} has the same cardinality as the set of all subsets of \mathbf{N} , thus $\text{card}(\mathbf{R}) = 2^{\aleph_0}$. The points in the Euclidean plane have the same cardinality: $\text{card}(\mathbf{R}^2) = \text{card}(\mathbf{R})$. But the set of all subsets of the plane has the larger cardinality $2^{2^{\aleph_0}}$, which is too much for our brains to keep track of. There are simply too many sets in the plane: we need to restrict attention to some not too large subclass of nice sets. For instance, the set of all disks has a much smaller cardinality, because three numbers suffice to determine a disk in the plane: its radius and the two coordinates of its center. This leads to the idea of simplifying a general, all too wild set, to some reasonable, more well-behaved set.

If we consider digital geometry we can make a similar discussion about cardinality. On a computer screen with, say, 1,024 pixels in a horizontal row and 768 pixels in a vertical column there are $1,024 \times 768 = 786,432$ pixels. To describe a subset of the screen we need to specify for each pixel whether it belongs to the set or not. This means that there are $2^{786,432} \approx 10^{236,740}$ different subsets of the screen, or binary images. A binary image is a black-and-white image, i.e., we only specify whether a point belongs to the set or to the background. (If we want to consider gray-level images, or color images, the cardinality goes up of course.) The really pure mathematician then says: there are only finitely many binary images. But the number of binary images must be compared with other finite numbers. Some astronomers estimate the mass of the universe at 10^{53} kg, which is 6×10^{79} proton masses or 10^{83} electron masses. One sometimes talks about “astronomical numbers,” which Webster defines as “enormously or inconceivably

large” numbers. This metaphor has not only faded; it is actually misleading in the world of image analysis.¹

Thus, although the number of binary images on a computer screen is finite, it is so huge that the conclusion must be the same as in the case of the infinite cardinal $2^{2^{80}}$: there are too many, and we must simplify. This leads, again, to image processing and mathematical morphology, this time of subsets of \mathbf{Z}^2 , or, a little more generally, of \mathbf{Z}^n , the set of all n -tuples of integers.

When we discuss mathematical morphology we want to keep both cases in mind, i.e., both the vector space \mathbf{R}^n of all n -tuples of real numbers and the digital space \mathbf{Z}^n . The latter covers the case of rectangular pixels, but, as already pointed out, also triangular and hexagonal pixels. What is common to \mathbf{R}^n and \mathbf{Z}^n is that they form an abelian group. Therefore we shall always start the discussion assuming that the space, called *image carrier*, is just an abelian group. When necessary we shall then specialize to \mathbf{R}^n and \mathbf{Z}^n . In many ways, \mathbf{R}^n will guide us in the more difficult study of \mathbf{Z}^n .

Serra (1982:6–15) lists “four principles of quantification.” These are about our ways to gather information about the external world. They apply also, but not exclusively, to image analysis. Let us discuss them briefly in mathematical terms.

Serra’s first principle is “compatibility under translation.” Mathematically speaking, this means that if we translate the object (i.e., move it some distance without rotating it) and then do something to it, then the result should be the same as if we perform the two operations in the other order. For a mapping, this simply means that $f(A + b) = f(A) + b$, which we may express as $f \circ T_b = T_b \circ f$, thus a kind of commutativity, writing T_b for the translation, $T_b(A) = A + b$. We say that f commutes with translations. However, when it comes to applications of this idea to images on the computer screen, we are in deep trouble. On a finite screen $\{x \in \mathbf{Z}^2; 0 \leq x_1 < 1,024, 0 \leq x_2 < 768\}$ almost nothing can commute with translations. To escape from this difficulty we introduce the ideal, infinite, computer screen \mathbf{Z}^2 . This principle is equally useful in \mathbf{R}^n and \mathbf{Z}^n .

Serra’s second principle is “compatibility under change of scale.” For a mapping this means that it commutes with homotheties (or dilatations), i.e., mappings of the form $x \mapsto \lambda x$. So $f(\lambda A) = \lambda f(A)$ for (say) positive λ . This is not problematic if we are in \mathbf{R}^n , but it certainly is if we are in \mathbf{Z}^n . This principle will therefore need to be suitably interpreted in \mathbf{Z}^n .

The third principle is that of “local knowledge.” This principle says that in order to know some bounded part of $f(A)$, we shall not need to know all of A , only some bounded part of A . Mathematically speaking: for every bounded set Y , there exists a bounded set Z such that $f(A \cap Z) \cap Y = f(A) \cap Y$. To know the result $f(A)$ in an arbitrary bounded set Y , we need not know all of A ; it is enough to know how A looks in some bounded set Z . The principle of local knowledge shall therefore not be understood in the topological sense: the key notion is boundedness.

Serra’s fourth principle of quantification is that of “semicontinuity.” It means that if a decreasing sequence (A_j) of closed sets tends to a limit A , thus $A = \bigcap A_j$, then $f(A_j)$ tends to $f(A)$. We shall later express this as a semicontinuity.... To do so we

¹In case you think that this conclusion depends on the universe having a rather low density, please calculate the mass of a fictitious universe with a radius of 14×10^9 light years $\approx 1.3 \times 10^{26}$ m and the density of a neutron star, say 10^{17} kg/m³. The conclusion is the same.

must of course define a topology.

2. Morphological operations on sets and functions

2.1. Infimal convolution

Let G be a semigroup, with the operation written as juxtaposition. This means that for any two elements $x, y \in G$ there is defined a product in G , written xy ; the product must be associative, i.e., $x(yz) = (xy)z$ for all $x, y, z \in G$. That's all. We do not assume that there is a neutral element; we do not assume that there are inverses; we do not assume commutativity. With this meager structure we may nevertheless define the *product* of two subsets A and B as

$$(2.1) \quad AB = \{xy; x \in A, y \in B\}.$$

If G happens to be commutative, it is customary to write the semigroup operation as addition, and we then speak of *Minkowski addition* of two subsets;

$$(2.2) \quad A + B = \{x + y; x \in A, y \in B\}$$

is called the *Minkowski sum* of A and B . However, for a while let us not assume commutativity.

In this way the family $\mathcal{P}(G)$ of all subsets of G becomes a semigroup. The empty set is a *zero* of this semigroup: $A\emptyset = \emptyset A = \emptyset$ for all $A \in \mathcal{P}(G)$. If G happens to have a neutral element 1 , then $\{1\}$ is a neutral element of $\mathcal{P}(G)$: $\{1\}A = A\{1\} = A$ for all $A \in \mathcal{P}(X)$. In the commutative case we write $A + \emptyset = \emptyset$ and $A + \{0\} = A$.

The product AB of two subsets can be put into the wider framework of *infimal convolution*. Given two functions $f, g: G \rightarrow [-\infty, +\infty]$ defined on the semigroup and with values in the extended real line $[-\infty, +\infty] = \mathbf{R} \cup \{-\infty, +\infty\}$, we define a new function $h = f \square g$, called the *infimal convolution* of f and g , as

$$(2.3) \quad (f \square g)(z) = h(z) = \inf_{x, y \in G} (f(x) \dagger g(y); xy = z), \quad z \in G.$$

The infimum is taken over all elements $x, y \in G$ such that their product is z , the argument of h . If there are no such elements, the infimum is plus infinity; that element is by definition the infimum over the empty set. There is a complication if f takes the value $+\infty$ at x and g takes the value $-\infty$ at y . We resolve this conflict by declaring that $+\infty$ shall win. So $s \dagger t$ is the usual sum if s and t are real numbers; if only one is infinite or both are infinite of the same sign, the sum takes that value; if s and t are infinite of opposite signs, we define the sum to be $+\infty$. In this way, this operation, called *upper addition*, becomes an upper semicontinuous mapping from $[-\infty, +\infty]^2$ into $[-\infty, +\infty]$. It is easy to check that it is associative. Similarly we define *lower addition*, $s \ddagger t = -((-\infty) \dagger (-t))$; here minus infinity wins. In my experience, using upper and lower addition is the most convenient method to calculate with the two infinities. This will be apparent, I hope, when we want to show that infimal convolution is associative.

For surveys of the properties of infimal convolution we refer to Moreau (1970), Rockafellar (1970), or Strömberg (1996).

The points where f or g takes the value $+\infty$ play no role in the formation of the infimum: the definition of upper addition guarantees this. Removing these points therefore yields an equivalent definition:

$$(2.4) \quad (f \square g)(z) = \inf_{x,y \in G} (f(x) + g(y); xy = z, f(x) < +\infty, g(y) < +\infty), \quad z \in G.$$

The *effective domain*, written $\text{dom } f$, of a function $f: X \rightarrow [-\infty, +\infty]$ defined on an arbitrary set X is the set where it is strictly less than plus infinity:

$$(2.5) \quad \text{dom } f = \{x \in X; f(x) < +\infty\}.$$

With this notation we can write (2.4) as

$$(2.6) \quad (f \square g)(z) = \inf_{\substack{x \in \text{dom } f \\ y \in \text{dom } g \\ xy = z}} (f(x) + g(y)), \quad z \in G.$$

Here there is no doubt about the meaning of the sum; the disadvantage with this approach is that we have to remember, each time we take an infimum, over which set the variables range.

Intuitively, plus infinity corresponds to vacuum and $-\infty$ to an infinitely dense neutron star. This may appear to be upside down. However, we should think of the density as $e^{-f(x)}$, and then of course $e^{-(+\infty)} = 0$. Infimal convolution is related to a kind of supremal convolution of the functions e^{-f}, e^{-g} , viz.

$$\sup_y [e^{-f(y)} e^{-g(x-y)}] = e^{-(f \square g)(x)}.$$

The supremum is often comparable to integration in \mathbf{R}^n , which means that we sometimes have a remarkably good approximation

$$e^{(f \square g)(x)} = \sup_{y \in \mathbf{R}^n} [e^{-f(y)} e^{-g(x-y)}] \approx \int_{\mathbf{R}^n} e^{-f(y)} e^{-g(x-y)} dy = (e^{-f} * e^{-g})(x), \quad x \in \mathbf{R}^n,$$

where the asterisk denotes usual convolution.

Proposition 2.1. *Infimal convolution is associative: $(f_1 \square f_2) \square f_3 = f_1 \square (f_2 \square f_3)$.*

Proof. We start calculating:

$$((f_1 \square f_2) \square f_3)(z) = \inf_{y, x_3} \left(\left[\inf_{\substack{x_1, x_2 \\ x_1 x_2 = y}} (f_1(x_1) \dot{+} f_2(x_2)) \right] \dot{+} f_3(x_3) \right).$$

According to the following lemma this expression is equal to

$$\inf_{\substack{y, x_3 \\ y x_3 = z}} \inf_{x_1, x_2} ([f_1(x_1) \dot{+} f_2(x_2)] \dot{+} f_3(x_3)) = \inf_{\substack{x_1, x_2, x_3 \\ x_1 x_2 x_3 = z}} (f_1(x_1) \dot{+} f_2(x_2) \dot{+} f_3(x_3)).$$

A similar calculation shows that the last expression is also equal to $(f_1 \square (f_2 \square f_3))(z)$. This proves associativity.

At a point in the proof we needed the following result.

Lemma 2.2. For any element $c \in [-\infty, +\infty]$ and any function $f: X \rightarrow [-\infty, +\infty]$ defined on an arbitrary set X we have

$$\inf_{x \in X} (c \dot{+} f(x)) = c \dot{+} \inf_{x \in X} f(x).$$

Proof. We just need to check all possibilities where our intuition is less reliable than usual, i.e., $c = \pm\infty$ or X empty.

Note that there are no exceptions to this formula. As a nasty little exercise, try to find the exact conditions under which the equality $\sup_x (c \dot{+} f(x)) = c \dot{+} \sup_x f(x)$ holds.

If G is a group, we know that $xy = z$ if and only if $x = zy^{-1}$, which in turn is equivalent to $y = x^{-1}z$, so the definition of $f \square g$ can be written

$$(f \square g)(z) = \inf_{x \in G} (f(x) \dot{+} g(x^{-1}z)) = \inf_{y \in G} (f(zy^{-1}) \dot{+} g(y)), \quad z \in G.$$

In the case of an abelian group this formula reads

$$(f \square g)(z) = \inf_{x \in G} (f(x) \dot{+} g(z - x)) = \inf_{y \in G} (f(z - y) \dot{+} g(y)), \quad z \in G.$$

Now why is infimal convolution more general than the product of sets or Minkowski addition? This is because of the formula

$$(2.7) \quad \text{dom}(f \square g) = \text{dom } f + \text{dom } g,$$

which is easily proved. A special case of this formula is obtained when we consider indicator functions.

To any subset A of a set X we define its *indicator function* i_A , which is simply defined as $i_A(x) = 0$ when $x \in A$ and $i_A(x) = +\infty$ when $x \notin A$. It is clear that $\text{dom } i_A = A$. We then have $i_A \square i_B = i_{AB}$ for all subsets A, B of a semigroup G , or, in the commutative case, $i_A \square i_B = i_{A+B}$. Hence the Minkowski sum may be defined in terms of infimal convolution as $A + B = \text{dom}(i_A \square i_B)$.

If G has a neutral element 1 , then $i_{\{1\}}$ is a neutral element for infimal convolution. More generally, if $f(x) = 0$ when $x \in A$ and $f(x) \geq 1$ otherwise and similarly with g and B , then $(f \square g)(z) = 0$ when $z \in AB$ and $(f \square g)(z) \geq 1$ when $z \notin AB$.

Thus multiplication of sets and Minkowski addition can be expressed in terms of infimal convolution. But we can go also in the other direction and express any infimal convolution as a product of sets, at the expense of adding one more dimension. This is done using the notion of strict epigraph.

Let us first define the *epigraph* of a function $f: X \rightarrow [-\infty, +\infty]$ defined on an arbitrary set X as

$$(2.8) \quad \text{epi } f = \{(x, t) \in X \times \mathbf{R}; f(x) \leq t\},$$

and the *strict epigraph* using instead strict inequality:

$$(2.9) \quad \text{epi}_s f = \{(x, t) \in X \times \mathbf{R}; f(x) < t\}.$$

If $X = G$ is a semigroup, then we make $G \times \mathbf{R}$ into a semigroup by defining $(x, s)(y, t) = (xy, s + t)$. It is not difficult to show that

$$\text{epi}_s(f \square g) = (\text{epi}_s f)(\text{epi}_s g).$$

This means that the function $f \square g$ can be defined as the function whose strict epigraph is the product $(\text{epi}_s f)(\text{epi}_s g)$. With the additive notation we have of course

$$\text{epi}_s(f \square g) = \text{epi}_s f + \text{epi}_s g.$$

For the epigraph we always have

$$\text{epi}(f \square g) \supset \text{epi} f + \text{epi} g,$$

where the inclusion relation may be strict (find examples).

2.2. Dilations and erosions

Let us now consider the commutative case and look at the Minkowski sum of two sets A and B :

$$(2.10) \quad A + B = \bigcup_{y \in B} (A + y) = \bigcup_{x \in A} (x + B).$$

Here $A + y = A + \{y\}$ is a *translate* of A : all points in A are moved by a fixed vector y . Thus the Minkowski sum is a union of translates of A , and equally a union of translates of B . We can view the sum in either way. We move A around in a manner determined by B or vice versa. Often A is some complicated set and B is a simple, small set, called the *structural element*. Then we say that $A + B$ is *A dilated* by B . Think of $G = \mathbf{R}^2$ or \mathbf{Z}^2 and with B equal to a ball. We fix B and apply the operation $A \mapsto A + B$ to a lot of sets A . We call this operation *dilation* by B .² Note that $A + \emptyset = \emptyset$ and that $A + \{0\} = A$. The mapping $A \mapsto A + B$ will be denoted by δ_B .³

We note that dilation by B commutes with translations:

$$T_y(A + B) = (A + B) + y = (A + y) + B = (T_y A) + B,$$

which we may write as $T_y \circ \delta_B = \delta_B \circ T_y$. It also commutes with the formation of unions:

$$(2.11) \quad \bigcup_{j \in J} (A_j + B) = \left(\bigcup_{j \in J} A_j \right) + B,$$

where $(A_j)_{j \in J}$ is a finite or infinite family of subsets of a semigroup G . Conversely, any mapping $f: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ which commutes with translations and the formation of infinite unions is dilation by $f(\{0\})$.

²This seems to be the most common term today; Matheron (1975:17) calls it *dilatation*, a term which is often reserved for mappings $x \mapsto \lambda x$.

³Sometimes $A + B$ is written $A \oplus B$; see Matheron (1975:16), Serra (1982, 2001) and Gonzalez & Woods (1993); note however that the sign \oplus is used in other contexts for the direct sum. There is no risk of misunderstanding $A + B$.

If B is finite, as is often the case in \mathbf{Z}^n , only finitely many checks are needed to decide whether a point x belongs to $A + B$: we check whether $x - b$ belongs to A for some $b \in B$. Therefore dilation by a finite set satisfies Serra's third principle (see the introduction, section 1.2).

To any given subset A of an abelian group G we define its *opposite set* as $\check{A} = \{-x; x \in A\}$. Then we can also define *Minkowski subtraction* by the formula

$$A - B = A + \check{B} = \{x - y; x \in A, y \in B\}.$$

Note that $\check{\check{B}} = \{0\} - B = 0 - B$; we may also write it as $-B$. We should not mix this up with the *set-theoretical minus*:

$$A \setminus B = A \cap \mathbf{C}B = A \cap B^c = \{x \in A; x \notin B\},$$

where we have written the complement of B in two different ways: $\mathbf{C}B = B^c$; both are quite usual. We now introduce another minus sign, written \ominus :

$$(2.12) \quad A \ominus B = \mathbf{C}(\mathbf{C}A + B) = (A^c + B)^c;$$

it is called the *erosion* of A by B . Note that $A \ominus \emptyset = G$ and that $A \ominus \{0\} = A$. The mapping $A \mapsto A \ominus B$ will be denoted by ε_B .

To illustrate the difference in meaning between these three minus signs, note that $A \setminus A = \emptyset$ for all sets A ; that $A - A \neq \emptyset$ if A is a nonempty subset of an abelian group; and that, finally, $A \ominus \check{A}$ is a subsemigroup if A is a subset of an abelian group G , cf. Lemma 2.3 below; if A is a bounded subset of \mathbf{R}^n or \mathbf{Z}^n for example, then $A \ominus \check{A} = \{0\}$.

Erosion is dual to dilation in a natural sense. Indeed, $A \ominus B = \mathbf{C}(\mathbf{C}A + B) = f^{\mathbf{d}}(A)$, where we define a mapping $f^{\mathbf{d}}$ for every mapping f by $f^{\mathbf{d}}(A) = \mathbf{C}f(\mathbf{C}A) = f(A^c)^c$; in this case $f(A) = A + B$ for a fixed B . Thus $\varepsilon_B = (\delta_B)^{\mathbf{d}} = \mathbf{C} \circ \delta_B \circ \mathbf{C}$. We shall say that $f^{\mathbf{d}}$ is the mapping *dual to* f .

We may also express an erosion as an intersection of translates of A . More precisely we have the following result.

Lemma 2.3. *For all subsets A and B of an abelian group G we have*

$$(2.13) \quad A \ominus B = \mathbf{C}(\mathbf{C}A + B) = \bigcap_{y \in B} (A + y) = \{x \in G; x + \check{B} \subset A\}.$$

Proof. Here the first equality is the definition of erosion. Let us call the third set Y . Then its complement is

$$\mathbf{C}Y = \mathbf{C} \bigcap_{y \in B} (A + y) = \bigcup_{y \in B} \mathbf{C}(A + y) = \bigcup_{y \in B} (\mathbf{C}A + y) = \mathbf{C}A + B = \mathbf{C}(A \ominus B),$$

so that $Y = A \ominus B$. Here we have used that translations commute with the operation of taking the complement: $\mathbf{C}(A + y) = \mathbf{C}A + y$.

Now $x \in Y$ means that $x \in A + y$ for all $y \in B$, which is equivalent to the statement that $x - y \in (A + y) - y = A$ for all $y \in B$. This, in turn, is equivalent to $x - B = x + \check{B} \subset A$, i.e., that x belongs to the last set in the statement of the lemma. Thus we have now proved the whole chain of equalities.

Using the duality of erosion and dilation it is easy to prove that erosion by B commutes with translations:

$$T_y(A \ominus B) = (T_y A) \ominus B,$$

succinctly, $T_y \circ \varepsilon_B = \varepsilon_B \circ T_y$. It also commutes with the formation of intersections:

$$(2.14) \quad \bigcap_{j \in J} (A_j \ominus B) = \left(\bigcap_{j \in J} A_j \right) \ominus B.$$

Conversely, any mapping $g: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ which commutes with translations and the formation of infinite intersections is erosion by $\mathbf{C}g(\mathbf{C}\{0\})$.

2.3. Exercises

2.1. Dilate and erode A by B when

(a) A is a disk in \mathbf{R}^2 , B a pair of points. (Here the word *disk* shall be understood with regard to one of the norms $\|\cdot\|_p$, $p = 1, 2, \infty$.)

(b) A is an open disk in \mathbf{R}^2 of radius R , B a closed disk of radius r .

(c) B is a translate of A .

(d) B is a translate of \check{A} .

2.2. (a) Prove that $A + B \subset C$ if and only if $A \subset C \ominus \check{B}$, where A, B, C are arbitrary subsets of an abelian group.

(b) Show by examples that it is not true that $A \ominus B \subset C$ if and only if $A \subset C + \check{B}$. Is one of the implications true?

(c) Give examples to show that we do not have an equivalence $A + B = C$ if and only if $A = C \ominus \check{B}$. However, the equivalence holds for certain families of sets...

2.3. Prove that, in an abelian group, erosion by B satisfies $\varepsilon_B(A) \subset A$ for all sets A if and only if $0 \in B$.

2.4. Calculate in \mathbf{R}^2 the Minkowski sum of two arbitrary segments $[a, b] = \{(1-t)a + tb; 0 \leq t \leq 1\}$ and $[c, d]$, where a, b, c, d are four arbitrary points in \mathbf{R}^2 .

2.5. Calculate in \mathbf{R}^2 the sum of three arbitrary segments.

2.6. Calculate in \mathbf{R}^3 the sum of three orthogonal segments.

2.7. Calculate in \mathbf{R}^3 the Minkowski sum of four segments, three of which are in one plane while the other is not.

2.8. (a) Calculate in \mathbf{R}^3 the sum of the four segments $[(0, 0, 0), (1, 0, 0)]$, $[(0, 0, 0), (0, 1, 0)]$, $[(0, 0, 0), (0, 0, 1)]$, and $[(0, 0, 0), (1, 1, 1)]$. The result is a polyhedron. How many vertices, edges, and faces does it have?

(b) Calculate in \mathbf{R}^3 the sum of the four unit normals to a regular tetrahedron.

2.9. (a) Calculate in \mathbf{R}^3 the Minkowski sum of the six segments $[(0, 0, 0), (1, 0, 0)]$, $[(0, 0, 0), (0, 1, 0)]$, $[(0, 0, 0), (0, 0, 1)]$, $[(0, 0, 0), (0, 1, 1)]$, $[(0, 0, 0), (1, 0, 1)]$, and $[(0, 0, 0), (1, 1, 0)]$. How many vertices, edges, and faces are there?

(b) Describe the Minkowski sum of the six edges of a regular tetrahedron in \mathbf{R}^3 .

2.10. Prove that a triangle cannot be the Minkowski sum of a finite number of segments.

2.11. A heptagon in \mathbf{R}^2 can never be the sum of a finite number of segments.

2.12. A tetrahedron cannot be the sum of a finite number of segments.

2.13. The vector sum of two triangles in \mathbf{R}^2 is a polygon. How many sides can it have? How many sides can the sum of a triangle and a square have?

2.14. Give an example of a quadrilateral in \mathbf{R}^2 that is not the sum of two segments.

3. Ordered sets: first notions

Let us define first the notion of preorder, which is even more fundamental than that of order. A *preorder* in a given set X is a relation (a subset of X^2) which satisfies two conditions: it is *reflexive* and *transitive*. This means, if we denote the relation by \leq , that for all $x, y, z \in X$,

$$(3.1) \quad x \leq x,$$

and

$$(3.2) \quad x \leq y \text{ and } y \leq z \text{ implies } x \leq z.$$

If, in addition, the relation is *antisymmetric*, i.e., satisfies

$$(3.3) \quad x \leq y \text{ and } y \leq x \text{ implies } x = y,$$

it is said to be an *order*. (Sometimes one says *partial order*.)

A *preordered set* is a set together with a preorder; an *ordered set* a set together with an order.

Suppose we have two preorders defined in a set X ; denote them by \leq and \preceq . The preorder \leq is said to be *finer* than the preorder \preceq , and \preceq is said to be *coarser* than \leq , if $x \leq y$ implies $x \preceq y$ for all x, y .

There is a finest preorder in a set, viz. when we define $x \leq y$ to mean that $x = y$. This preorder is of course an order; let us call it the *discrete order*. There is also a coarsest preorder in any set X , when we declare that $x \leq y$ for all $x, y \in X$. Let us call this the *chaotic preorder*.

Given any preorder we define $[a, \rightarrow[$ as the set of all x such that $a \leq x$ and $] \leftarrow, a]$ as the set of all x such that $x \leq a$. The sets $[a, \rightarrow[$ and $] \leftarrow, a]$ determine a for all $a \in X$ if and only if \leq is an order.

Given a preorder \leq in X , we can introduce an equivalence relation \simeq in X (a reflexive, symmetric and transitive relation) by defining $x \simeq y$ to mean that $x \leq y$ and $y \leq x$. If \leq is an order, then \simeq is just equality. We can form the quotient space X/\simeq of all equivalence classes of X modulo \simeq . The equivalent classes are just

$$[a, \rightarrow[\cap] \leftarrow, a] = \{x \in X; a \leq x \leq a\}, \quad a \in X.$$

In this quotient space, \leq induces an order.

If two preordered sets X, Y are given, we may consider mappings $f: X \rightarrow Y$, and of these the *increasing* (order preserving, isotone⁴) mappings are of interest; these are the mappings which satisfy

$$(3.4) \quad \text{for all } x, x' \in X, x \leq_X x' \text{ implies } f(x) \leq_Y f(x').$$

⁴Birkhoff 1948:49

We shall write $\text{Incr}(X, Y)$ for the set of all increasing mappings $X \rightarrow Y$. When $X = Y$ we shall sometimes write X' for $\text{Incr}(X, X)$. We may think of $\text{Incr}(X, Y)$ as an analogue of the linear mappings from a vector space into another.

A preorder \leq is finer than another preorder \preceq if and only if the identity $(X, \leq) \rightarrow (X, \preceq)$ is increasing.

Note that if X has the discrete order, where $x \leq y$ means $x = y$, then $\text{Incr}(X, Y)$ consists of all mappings $X \rightarrow Y$; $\text{Incr}(X, Y) = Y^X$. The conclusion is the same if Y is equipped with the chaotic preorder. If on the other hand Y has the discrete order, then $\text{Incr}(X, Y)$ consists of those mappings $X \rightarrow Y$ that are constant on any chain of comparable elements, thus $f(x_j)$ is constant if $x_1 \leq x_2 \geq x_3 \leq \dots \leq x_{2n}$. Let us say that a preordered set X is *connected* if, given any elements $a, b \in X$, there exists a finite chain of elements $x_1 \leq x_2 \geq x_3 \leq \dots \leq x_{2n}$ passing through a and b . Then any mapping $f \in \text{Incr}(X, Y)$ is constant if X is connected and Y has the discrete order. The same conclusion holds if X has the chaotic preorder and Y is ordered.

Another interesting property is obtained when we turn around the implication sign in (3.4). Let us agree to call a mapping $f: X \rightarrow Y$ *coincreasing* if it satisfies

$$(3.5) \quad \text{for all } x, x' \in X, f(x) \leq_Y f(x') \text{ implies } x \leq_X x'.$$

Increasing does not imply coincreasing; coincreasing does not imply increasing; cf. exercise 2.1. If X is ordered, then a mapping is coincreasing if and only if it is injective and the mapping $f(x) \mapsto x$ is increasing.

A basic example of an ordered set is the set $\mathcal{P}(W)$ of all subsets of a set W , with the order relation given by inclusion, thus $A \leq B$ being defined as $A \subset B$ for $A, B \in \mathcal{P}(W)$.

Definition 3.1. A *closure operator*⁵ (or *closing*⁶) in an ordered set X is a mapping $X \ni x \mapsto \bar{x} \in X$ which is *extensive*⁷ (or *expanding*), *increasing*, and *idempotent*; in other words, which satisfies the following three conditions for all $x, y \in X$:

$$(3.6) \quad x \leq \bar{x};$$

$$(3.7) \quad x \leq y \text{ implies } \bar{x} \leq \bar{y};$$

$$(3.8) \quad \overline{\bar{x}} = \bar{x}.$$

When checking (3.8) it is of course enough to prove that $\overline{\bar{x}} \leq \bar{x}$ if we have already proved (3.6).

The element \bar{x} is said to be the *closure* of x . Elements x such that $\bar{x} = x$ are called *closed* (for this operator). An element is closed if and only if it is the closure of some element (and then it is the closure of itself).

⁵Ore (1944:494) and Birkhoff (1948:49) used the term *closure operation*; the latter attributed the concept to Moore (1910:53–80). However, although Moore wrote about *closure* and *closure properties*, he did not give a clearcut definition. Everett (1944) used *closure operator*. Dubreil & Dubreil-Jacotin (1964:9, 177) calls the operator *fermeture de Moore*.

⁶This term seems to have been introduced by Matheron (1975:18). In Matheron (1967:18) he used the French term *fermeture* for the notion.

⁷Birkhoff (1948:49).

Sometimes we shall consider operators that are defined only on a subset of X .

Let $f: Y \rightarrow X$ be a mapping, where $Y \subset X$. Then it is not obvious what idempotency shall mean, for $f(f(y))$ need not be defined. We solve this problem by a modified definition; see Kiselman (1969:336).⁸

Definition 3.2. *We shall say that $f: Y \rightarrow X$, where Y is a subset of an ordered set X , is a closure operator if it satisfies the following two conditions for all $x, y \in Y$:*

$$(3.9) \quad x \leq f(x);$$

$$(3.10) \quad x \leq f(y) \text{ implies } f(x) \leq f(y).$$

It follows from this that f is increasing, and that $f(f(x)) = x$ for every $y \in Y$ such that $f(y)$ happens to belong to Y . It is also clear that if $Y = X$, then the new definition agrees with Definition 3.1. Any closure operator f in the sense of Definition 3.2 can be extended to an idempotent operator $f_1: Z \rightarrow X$ by taking $Z = Y \cup \text{im} f$ and defining $f_1(x) = x$ when $x \in Z \setminus Y$.

A basic example of a closure operator is the topological closure operator which associates to a set in a topological space its topological closure, i.e., the smallest closed set containing the given set. In fact a closure operator in $\mathcal{P}(W)$ defines a topology in W if and only if it satisfies, in addition to (3.6), (3.7), (3.8) above, two extra conditions, viz. that $\overline{\emptyset} = \emptyset$ and

$$\overline{A \cup B} = \overline{A} \cup \overline{B} \text{ for all } A, B \subset W.$$

Another closure operator of great importance is the operator which associates to a set A in \mathbf{R}^n its convex hull, the smallest convex set containing the given set, denoted by $\text{cvx}A$. The composition $A \mapsto \overline{\text{cvx}A}$ is a closure operator, whereas the composition in the other order, $A \mapsto \text{cvx}(\overline{A})$ is not idempotent (find examples). We see that the composition of two closure operators is sometimes, but not always, a closure operator.

In both these examples X is the power set of some set W , and the closure operator is given as an intersection:

$$\overline{A} = \bigcap \{Y; Y \text{ is closed and } Y \supset A\}.$$

More generally, if a closure operator $x \mapsto \bar{x}$ is given and we denote by F the set of its closed elements, then

$$(3.11) \quad \bar{x} = \inf\{y \in F; y \geq x\}.$$

Conversely, any subset F of X such that the infimum⁹ of a subset of F always exists defines a closure operator by formula (3.11).

Dual to the notion of closure operator is the notion of *opening*.

Definition 3.3. *A mapping $X \ni x \mapsto x^\circ \in X$ is said to be an opening¹⁰ if it is antiextensive, increasing, and idempotent; in other words, if it satisfies the following*

⁸Matheron (1975:186) required instead that $f(y)$ belong to Y whenever $y \in Y$.

⁹For the definition of infimum in an ordered set, see Chapter 7

¹⁰This term seems to have been introduced by Matheron (1975:18). In his earlier book (1967:18) he used the French term *ouverture*. Tucker (1936:94) used the term *aperture* for the dual of closure. However, I do not know how to name the operator $x \mapsto x^\circ$ using this word.

three conditions for all $x, y \in X$:

$$(3.12) \quad x^\circ \leq x;$$

$$(3.13) \quad x \leq y \text{ implies } x^\circ \leq y^\circ;$$

$$(3.14) \quad (x^\circ)^\circ = x^\circ.$$

Definition 3.4. Let $g: Y \rightarrow X$ be a mapping of a subset Y into X . We shall say that g is an opening if it satisfies the following two conditions for all $x, y \in Y$:

$$(3.15) \quad g(x) \leq x;$$

$$(3.16) \quad g(x) \leq y \text{ implies } g(x) \leq g(y).$$

4. Closure operators on subsets of an abelian group

4.1. Combining erosions and dilations

We have defined in section 2 dilations and erosions. We shall now combine them.

It is clear that $(A + B) + C = A + (B + C)$, so the composition of two dilations is a dilation: $\delta_C \circ \delta_B = \delta_{B+C}$. Similarly, $(A \ominus B) \ominus C = A \ominus (B + C)$, so the composition of two erosions is an erosion: $\varepsilon_C \circ \varepsilon_B = \varepsilon_{B+C}$. But what about the composition of a dilation and an erosion? Let us consider $f = \varepsilon_C \circ \delta_{\check{B}}$, defined by $f(A) = (A - B) \ominus C$, where B and C are fixed subsets of an abelian group. This operation is always increasing, and it is extensive if and only if $C \subset B$, which in turn is equivalent to $f(\{0\}) \supset \{0\}$. We shall now study this operator when $C = B$ and see that it is idempotent then (cf. Exercise 4.4). We shall also see that, in typical cases, it is neither a dilation nor an erosion.

Given two subsets A and B of an abelian group G we define

$$(4.1) \quad A^B = (A + \check{B}) \ominus B = (A - B) \ominus B = \mathbf{C}(\mathbf{C}(A - B) + B) = \{x \in G; x - B \subset A - B\}.$$

This means that we perform on A first a dilation by \check{B} , then an erosion by B :

$$(4.2) \quad A^B = \varepsilon_B(\delta_{\check{B}}(A)) = (\varepsilon_B \circ \delta_{\check{B}})(A).$$

(In many applications, B is symmetric, so that $\check{B} = B$.) We shall see that this is a closure operator in the sense of section 3. So we call A^B the B -closure of A , and a set A is called B -closed if $A^B = A$.¹¹

Analogously we define

$$(4.3) \quad A_B = (A \ominus \check{B}) + B = \mathbf{C}(\mathbf{C}A - B) + B = \bigcup_{x \in G} (x + B; x + B \subset A)$$

¹¹Matheron (1975) and Serra (1982, 2001) use the notation A^B ; Gonzalez & Woods (1993:524) write $A \bullet B$.

by performing first an erosion by \check{B} and then a dilation by B :

$$(4.4) \quad A_B = \delta_B(\varepsilon_{\check{B}}(A)) = (\delta_B \circ \varepsilon_{\check{B}})(A).$$

We call A_B the B -opening of A , and a set A is called B -open if $A_B = A$.¹²

Since ε_B is the mapping dual to δ_B , i.e., $\varepsilon_B = (\delta_B)^\mathfrak{d} = \mathfrak{C} \circ \delta_B \circ \mathfrak{C}$, we have

$$\delta_B \circ \varepsilon_{\check{B}} = \delta_B \circ \mathfrak{C} \circ \delta_{\check{B}} \circ \mathfrak{C} = \mathfrak{C} \circ \varepsilon_B \circ \delta_{\check{B}} \circ \mathfrak{C},$$

which implies that $\mathfrak{C}(A_B) = (\mathfrak{C}A)^B$ and shows that a set is B -open if and only if its complement is B -closed.

Theorem 4.1. *For any subset B of an abelian group G , the mapping $\mathcal{P}(G) \ni A \mapsto A^B \in \mathcal{P}(G)$ is a closure operator, and the mapping $A \mapsto A_B$ is an opening.*

The theorem will follow on combining Propositions 4.2 and 4.4 below.

Proposition 4.2. *For any subset B of an abelian group G , the mappings $A \mapsto A^B$ and $A \mapsto A_B$ are increasing; the first is extensive and the second is antiextensive.*

Proof. It is obvious that the mappings are increasing. The last formula in (4.3) shows that A_B is contained in A . That $A \subset A^B$ follows by duality.

Proposition 4.3. *For any subsets A, B of an abelian group G , $A + B$ is A -open and B -open, and $A \ominus B$ is $\mathfrak{C}A$ -closed and B -closed.*

Proof. We always have $(A + B)_B \subset A + B$ (Proposition 4.2). On the other hand, $(A + B)_B = ((A + B) \ominus \check{B}) + B = A^{\check{B}} + B \supset A + B$ by the same proposition, so we have equality: $(A + B)_B = A + B$, meaning that $A + B$ is B -open. By symmetry it is also A -open.

Passing to complements, we deduce that $\mathfrak{C}(A + B)$ is A -closed and B -closed. But $\mathfrak{C}(A + B) = \mathfrak{C}A \ominus B$, so the set $\mathfrak{C}A \ominus B$ is A -closed and B -closed. We now just replace A by $\mathfrak{C}A$ to get the second statement of the proposition.

Proposition 4.4. *For any subsets B of an abelian group G , the mappings $A \mapsto A^B$ and $A \mapsto A_B$ are idempotent.*

Proof. By Proposition 4.3, $A^B = (A - B) \ominus B$ is B -closed, which means that $(A^B)^B = A^B$. The statement about A_B follows by duality.

For any fixed B , the family of all sets $A \ominus B$ is equal to the family of all B -closed sets. This is the same as the family of all sets A^B . Suppose we put on spectacles blocking out everything except the B -closed sets. The whole world becomes B -closed: our B -spectacles do not permit us to see anything else. Then we can only see sets of the form $A \ominus B$. However, A^B is also B -closed and usually a much better approximation of A than $A \ominus B$.

¹²Matheron (1975) and Serra (1982, 2001) use the notation A_B , Gonzalez & Woods (1993:524) use $A \circ B$. Note, however, that \circ is already used to denote an operation in mathematics, viz. composition of functions or relations.

4.2. Characterizing closure operators which commute with translations

In an abelian group the mappings that commute with translations are of special significance. We can characterize closure operators which commute with translations in terms of Minkowski addition as follows.

Theorem 4.5. *Let G be an abelian group and \mathcal{A} a subfamily of $\mathcal{P}(G)$ such that \mathcal{A} contains all singleton sets $\{x\}$, $x \in G$, and such that $A + B \in \mathcal{A}$ for all $A, B \in \mathcal{A}$. Let $f: \mathcal{A} \rightarrow \mathcal{P}(G)$ be a closure operator in the sense of Definition 3.2. Then the following conditions are equivalent.*

- (A) f commutes with all translations;
- (B) $f(A + x) \supset f(A) + x$ for all $A \in \mathcal{A}$ and all $x \in G$;
- (C) $f(A + B) \supset f(A) + B$ for all $A, B \in \mathcal{A}$ (thus $f \circ \delta_B \supset \delta_B \circ f$ for all $B \in \mathcal{A}$);
- (D) $f(A + B) \supset f(A) + f(B)$ for all $A, B \in \mathcal{A}$.

Proof. That (A) and (B) are equivalent is easily proved.

Assume that (B) holds. Then $f(A + B) \supset f(A + y) \supset f(A) + y$ for all $y \in B$. Now (C) follows on taking the union over all y .

If (C) holds we know that $f(A + B) \supset A + f(B)$, which implies that $f(A + B) \supset f(A + f(B))$ by (3.8). Applying (C) a second time we see that the latter set contains $f(A) + f(B)$; hence (D) holds.

Finally, if (D) holds, then $f(A + x) \supset f(A) + g(\{x\}) \supset f(A) + x$. Thus (B) holds.

In the case of a vector space and a mapping which commutes with homotheties, thus $f(\lambda A) = \lambda f(A)$ for nonnegative λ , it follows from (D) that $f(\lambda A + \mu B) \supset \lambda f(A) + \mu f(B)$ for all $\lambda, \mu \geq 0$, which is a property analogous to concavity.¹³

4.3. Matheron's structure theorems

We are now acquainted with dilations, erosions, closings, and openings. They are examples of morphological mappings, but we have so far no idea how general they are. Are they just some special mappings that we have come across? There are obviously many more mappings $\mathcal{P}(G) \rightarrow \mathcal{P}(G)$. It is reasonable to limit our consideration to increasing mappings which commute with all translations. How special are the dilations and erosions etc. we have studied so far in the family of these more general mappings?

Georges Matheron found a structural theorem which describes the structure of increasing mappings in \mathbf{R}^n which commute with translations in terms of dilations and erosions. Actually the result holds in every abelian group.

For any mapping $f: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ of the family of all subsets of an abelian group G into itself we define its *kernel* as $\ker f = \{A \in \mathcal{P}(G); 0 \in f(A)\}$. The notion is due to Matheron (1975:217). A mapping which commutes with all translations is completely determined by its kernel.

We mention some examples of kernels.

1. The kernel of the identity $\text{Id}: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ is $\ker \text{Id} = \{A \in \mathcal{P}(G); 0 \in A\}$.
2. The kernel of a translation T_b , $T_b(A) = A + b$, is $\ker T_b = \{A; b \in A\}$.
3. A dilation $\delta_{\check{B}}: A \mapsto A + \check{B}$ has kernel $\ker(\delta_{\check{B}}) = \{A; A \cap B \neq \emptyset\}$.
4. The kernel of an erosion $\varepsilon_{\check{B}}: A \mapsto A \ominus \check{B}$ is $\ker(\varepsilon_{\check{B}}) = \{A; A \supset B\}$.

¹³This property was introduced and used to study closure operators appearing in the theory of partial differential equations in Kiselman (1969); however, I did not know at the time that it is implied by (A).

5. The kernel of a closure operator $f = \varepsilon_B \circ \delta_{\check{B}}$ is $\ker f = \bigcap_{y \in B} \ker(\delta_{B-y})$.
6. An opening $g = \delta_B \circ \varepsilon_{\check{B}}: A \mapsto A_B$ has kernel $\ker g = \bigcup_{y \in B} \ker(\varepsilon_{B-y})$.

Matheron's first structural theorem (below) shows that very general mappings can be expressed in terms of the dilations and erosions we have already studied. It is important because it underlines the fact that the mappings we already know are the building blocks of a much more general class of mappings.

Theorem 4.6 (Matheron (1975: Proposition 8-1-3)). *Let G be an abelian group and $f: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ an increasing mapping which commutes with translations. Then f is a union of erosions as well as an intersection of dilations:*

$$(4.5) \quad f(A) = \bigcup_{B \in \ker f} (A \ominus \check{B}) = \bigcap_{B \in \ker f^{\mathbf{d}}} (A - B), \quad A \in \mathcal{P}(G),$$

where $f^{\mathbf{d}}(A) = \mathbb{C}f(\mathbb{C}A)$ is the mapping dual to f .

Proof. If $B \in \ker f$, then $0 \in f(B)$ and, since f commutes with all translations, $x \in f(x + B)$ for all $x \in G$. Therefore, given any subset A of G , $x \in A \ominus \check{B}$ implies that $x + B \subset A$ and $x \in f(x + B) \subset f(A)$ (f is increasing). Since x is arbitrary in $A \ominus \check{B}$, we have proved that $A \ominus \check{B} \subset f(A)$ for all A and all $B \in \ker f$. Letting B vary in $\ker f$ we see that the union of all the $A \ominus \check{B}$ is contained in $f(A)$.

To prove the inclusion in the other direction, take an arbitrary element x of $f(A)$. Then $0 \in f(A - x)$, i.e., $A - x \in \ker f$. If we now define B as $A - x$, then $x \in A \ominus \check{B}$. This means that there exists a $B \in \ker f$ such that $x \in A \ominus \check{B}$, and we have proved that $f(A)$ is contained in the union of all the $A \ominus \check{B}$.

The second representation follows on applying the first to $f^{\mathbf{d}}$.

Let X be an arbitrary set and \mathcal{A} a subset of $\mathcal{P}(X)$. For any mapping $f: \mathcal{A} \rightarrow \mathcal{P}(X)$ we then define

$$(4.6) \quad f_{\diamond}(B) = \bigcup_{\substack{A \in \mathcal{A} \\ A \subset B}} f(A) \quad \text{and} \quad f^{\diamond}(B) = \bigcap_{\substack{A \in \mathcal{A} \\ A \supset B}} f(A).$$

Lemma 4.7. *The mappings f_{\diamond} and f^{\diamond} are dual to each other in the sense that $(f_{\diamond})^{\mathbf{d}} = (f^{\mathbf{d}})^{\diamond}$ and $(f^{\diamond})^{\mathbf{d}} = (f^{\mathbf{d}})_{\diamond}$. If $f: \mathcal{A} \rightarrow \mathcal{P}(X)$ is an increasing mapping, then f_{\diamond} is the smallest increasing extension of f to all of $\mathcal{P}(X)$, and f^{\diamond} is the largest increasing extension of f to all of $\mathcal{P}(X)$.*

Proof. To prove the duality result is straightforward. If f is increasing, it follows that f_{\diamond} is actually an extension of f , i.e., that $f_{\diamond}|_{\mathcal{A}} = f$. If g is an arbitrary increasing extension, we see that we must have $g(B) \supset f_{\diamond}(B)$ for all B , proving that f_{\diamond} is the smallest increasing extension. The result for f^{\diamond} follows by duality.

Proposition 4.8. *Let X be any set, \mathcal{A} a subset of $\mathcal{P}(X)$, and let $g: \mathcal{A} \rightarrow \mathcal{P}(X)$ be an opening in the sense of Definition 3.4. Assume that $g(A) \in \mathcal{A}$ for every $A \in \mathcal{A}$. Then $g_{\diamond}: \mathcal{P}(X) \rightarrow \mathcal{P}(X)$ is an opening.*

Proof. That g_{\diamond} is increasing and antiextensive is obvious. It remains to be proved that g_{\diamond} is idempotent. We know that g_{\diamond} is an extension of g , which implies that

$g_\diamond(g(A)) = g(g(A)) = g(A)$ for all $A \in \mathcal{A}$, since by hypothesis $g(A) \in \mathcal{A}$. Hence $g_\diamond(g_\diamond(B)) \supset g_\diamond(g(A)) = g(A)$ if $A \in \mathcal{A}$ and $A \subset B$. On taking the union over all $A \in \mathcal{A}$ contained in B we get

$$g_\diamond(g_\diamond(B)) \supset \bigcup_{\substack{A \in \mathcal{A} \\ A \subset B}} g(A) = g_\diamond(B) \supset g_\diamond(g_\diamond(B)),$$

showing that g_\diamond is idempotent.

Remark. The hypothesis that g maps \mathcal{A} into itself cannot be dispensed with. Indeed, take $X = [0, +\infty[$, $A_j = [0, j]$ for $j = 0, 1, 2, 3$, and define $g(A_j) = A_{j-1} = [0, j-1]$, $j = 1, 3$, and $\mathcal{A} = \{A_1, A_3\}$. Then $g: \mathcal{A} \rightarrow \mathcal{P}(X)$ is an opening in the sense of Definition 3.4; in particular (3.14) holds: $g(A_j) \subset A_k$ implies $g(A_j) \subset g(A_k)$ for $j, k = 1, 3$. But g_\diamond is not idempotent: $A_0 = g_\diamond(g_\diamond(A_3)) \neq g_\diamond(A_3) = A_2$.

For any mapping $f: \mathcal{A} \rightarrow \mathcal{P}(G)$ we shall say that a set A is *f-invariant* if $f(A) = A$, and we shall denote by $\text{Inv}_f = \{A \in \mathcal{A}; f(A) = A\}$ the set of all *f-invariant* sets. If f is a closing, this is the set of all *f-closed* elements; if f is an opening, it is the set of all *f-open* elements.

Proposition 4.9. *If $g: \mathcal{A} \rightarrow \mathcal{P}(X)$ is an opening, then Inv_g is closed under the formation of unions (in particular $\emptyset \in \text{Inv}_g$), and $g = (\text{Id}_{\text{Inv}_g})_\diamond \big|_{\mathcal{A}}$, i.e., g is the smallest increasing extension to \mathcal{A} of the identity on Inv_g . Conversely, let a class $\mathcal{C} \subset \mathcal{P}(X)$ be given. Then $(\text{Id}_{\mathcal{C}})_\diamond$, the smallest increasing extension of the identity on \mathcal{C} , is an opening h , and Inv_h is the class closed under the formation of unions generated by \mathcal{C} .*

Proof. If $A_j \in \mathcal{A}$, $j \in J$, are g -invariant, then $g(\bigcup A_j) \supset g(A_k)$ for all k , so $g(\bigcup A_j) \supset \bigcup g(A_k) = \bigcup A_k$. Since on the other hand we always have $g(\bigcup A_j) \subset \bigcup A_j$, it follows that $\bigcup A_j \in \text{Inv}_g$. (We have actually only used (3.12) and (3.13) here.)

That g is the smallest increasing extension of the identity on Inv_g follows from

$$(\text{Id}_{\text{Inv}_g})_\diamond(B) = \bigcup_{\substack{A \in \mathcal{A} \\ A \subset B}} A = B \text{ if } B \in \mathcal{A}.$$

Since the inclusion mapping $\mathcal{C} \rightarrow \mathcal{P}(X)$ is an opening, it follows from Proposition 4.7 that $(\text{Id}_{\mathcal{C}})_\diamond$ is an opening. To prove the last assertion in the statement of the proposition, we note that by the definition of $(\text{Id}_{\mathcal{C}})_\diamond$, $h(B) = B$ implies that $B = h(B) = \bigcup(A; A \in \mathcal{C}, A \subset B)$, which means that B belongs to the set closed under union formation generated by \mathcal{C} .

In case $X = G$ is an abelian group, let us say that a subset \mathcal{B} of Inv_g is a *basis for Inv_g* if Inv_g is the class closed under formation of unions and translations generated by \mathcal{B} . This means that Inv_g is the family of all sets of the form $\bigcup_{(B,x) \in M} (B+x)$ for some subset M of $\mathcal{B} \times G$.

Matheron also proved a second structure theorem, which characterizes general (algebraic) closings and openings in terms of the elementary operations $A \mapsto A^B, A_B$, the closings and openings by a given subset B . His result, which shows that the latter are not so special as one could imagine, is as follows.

Theorem 4.10 (Matheron (1975: Proposition 7-1-3]). *Let G be an abelian group and $g: \mathcal{P}(G) \rightarrow \mathcal{P}(G)$ a mapping which commutes with translations. Then g is an opening if and only if it admits a representation*

$$g(A) = \bigcup_{B \in \mathcal{B}} A_B$$

for some class $\mathcal{B} \subset \mathcal{P}(G)$. If this is so, the class of sets invariant under g is the family of all elements of \mathcal{B} and all translates and unions of these. If g is an opening, the dual mapping g^d is a closing and has the representation

$$g^d(A) = \bigcap_{B \in \mathcal{B}} A^B.$$

Proof. First assume that g is an opening. The smallest extension of the identity on Inv_g is given by

$$(4.7) \quad g(A) = (\text{Id}_{\text{Inv}_g})_{\diamond}(A) = \bigcup_{\substack{C \in \text{Inv}_g \\ C \subset A}} C = \bigcup_{B \in \mathcal{B}} \bigcup_{x \in G} (B + x; B + x \subset A),$$

if \mathcal{B} is a basis for Inv_g as defined before the statement of the theorem. Let us now note that

$$A_B = (A \ominus \check{B}) + B = \bigcup_{x \in G} (B + x; B + x \subset A).$$

Therefore the last expression in (4.7) is equal to $\bigcup_{B \in \mathcal{B}} A_B$.

Conversely, if $g_j, j \in J$, are openings, so is g defined by $g(A) = \bigcup g_j(A)$. Indeed, g is certainly increasing and antiextensive, and idempotency follows from

$$g(g(A)) = \bigcup_j g_j \left(\bigcup_k g_k(A) \right) \supset \bigcup_j g_j(g_j(A)) = \bigcup_j g_j(A) = g(A).$$

4.4. Exercises

4.1. Define in an abelian group G a mapping $g = \delta_C \circ \varepsilon_{\check{B}}$, thus $g(A) = (A \ominus B) + C$, $A \in \mathcal{P}(G)$, where B and C are fixed subsets of G . Prove that g is antiextensive if and only if $C \subset B$, which in turn is equivalent to $g(G \setminus \{0\}) \subset G \setminus \{0\}$. Prove that if B and C are nonempty bounded subsets of \mathbf{Z}^n or \mathbf{R}^n with $C \subset B$, then g is idempotent only if $B = C$, thus only if $g(A) = A_B$. *Hint:* Calculate $g(B)$ and $g(g(B))$.

4.2. Study the dual mapping $f = \varepsilon_C \circ \delta_{\check{B}}$.

5. Distance transforms

5.1. Definition and basic properties of distance transforms

Distance transformations of digital images are a useful tool in image analysis. The distance transform of a set (or shape, or image) is a function on the image carrier. Outside the set, the value of the distance transform at a certain pixel is defined to be

the distance from that pixel to the set. Inside the set, it is often defined as the distance to the complement, but we shall find it convenient to define it as minus the distance to the complement, for reasons that will become clear soon.

The distances can be measured in different ways, e.g., by approximating the Euclidean distance in the two-dimensional image, the Euclidean distance between two pixels $x = (x_1, x_2)$ and $y = (y_1, y_2)$ being

$$\|x - y\|_2 = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}.$$

Other distances that have been used are the *city-block distance* or l^1 -distance

$$\|x - y\|_1 = |x_1 - y_1| + |x_2 - y_2|$$

and the *chessboard distance* or l^∞ -distance

$$\|x - y\|_\infty = \max(|x_1 - y_1|, |x_2 - y_2|).$$

Let X be any nonempty set. Let us agree to call a function $d: X \times X \rightarrow \mathbf{R}$ a *distance* if d is positive definite:

$$(5.1) \quad d(x, y) \geq 0 \text{ with equality precisely when } x = y, \quad x, y \in X,$$

and symmetric:

$$(5.2) \quad d(x, y) = d(y, x) \text{ for all } x, y \in X.$$

A distance will be called a *metric* if in addition it satisfies the triangle inequality:

$$(5.3) \quad d(x, z) \leq d(x, y) + d(y, z) \text{ for all } x, y, z \in X.$$

Every nonempty set can be equipped with a metric, viz. the discrete metric d_0 defined as

$$(5.4) \quad d_0(x, x) = 0, \quad d_0(x, y) = 1 \text{ if } x \neq y.$$

We shall use the word *distance* in a noncommittal way, thus not implying that it is a metric.

The set X will usually be the image plane \mathbf{Z}^2 consisting of all points in the plane with integer coordinates (the pixels), or more generally the image space \mathbf{Z}^n . Whenever X is an abelian group it is of particular interest to use translation-invariant distances, i.e., those which satisfy

$$(5.5) \quad d(x - a, y - a) = d(x, y) \text{ for all } a, x, y \in X.$$

A *metric space* is simply a set provided with a metric. In any metric space we define balls as follows: the *closed ball of center c and radius r* is the set of all points x satisfying $d(x, c) \leq r$ and will be denoted by $B_{\leq}(c, r)$; the *open ball of center c and radius r* is the set of all points satisfying the strict inequality $d(x, c) < r$; it will be denoted by $B_{<}(c, r)$. In the Euclidean case these notions are well known, but in general we must be a bit careful: the closure of $B_{<}(c, r)$ with respect to the topology defined by

d is not necessarily equal to $B_{\leq}(c, r)$, and the interior of $B_{\leq}(c, r)$ is not necessarily equal to $B_{<}(c, r)$ (construct examples). Also note that if two balls $B_{<}(c_1, r_1)$ and $B_{<}(c_2, r_2)$ with $r_1, r_2 > 0$ are disjoint, then we can only conclude that $\max(r_1, r_2) \leq d(c_1, c_2)$, whereas in a normed space a stronger inequality, $\max(r_1, r_2) \leq r_1 + r_2 \leq \|c_1 - c_2\|$, holds.

Every metric defines a topology: a set is declared to be open if and only if it is a union of open balls. However, we shall sometimes use another topology on X than that defined by d .

We note that in any abelian group with a translation-invariant metric we have the relations

$$B_{<}(c_1, r_1) + B_{<}(c_2, r_2) \subset B_{\leq}(c_1, r_1) + B_{<}(c_2, r_2) \subset B_{<}(c_3, r_3);$$

$$B_{\leq}(c_1, r_1) + B_{\leq}(c_2, r_2) \subset B_{\leq}(c_3, r_3),$$

where $c_1 + c_2 = c_3$, $r_1 + r_2 = r_3$. In a vector space, with d defined by $d(x, y) = \|x - y\|$ using some norm $\|\cdot\|$, the inclusions here are actually equalities if $r_1, r_2 > 0$.

If a metric is defined we define the *distance transform* DT_A of a subset A of X by

$$(5.6) \quad \text{DT}_A(x) = \begin{cases} -\inf_{y \notin A} d(x, y) = -\sup(r; B_{<}(x, r) \subset A), & x \in A; \\ \inf_{y \in A} d(x, y) = \sup(r; B_{<}(x, r) \subset \complement A), & x \in X \setminus A. \end{cases}$$

Thus when x does not belong to a set A , $\text{DT}_A(x)$ is what we naturally understand by the distance from x to A , but we complement this idea by defining also the distance transform inside A . It is then natural to take it negative there, so that A is approximately the set where the transform is negative. It is also very convenient that, in a normed space, the distance transform of a convex set is a convex function; see Corollary 5.8.

Note the symmetry: $\text{DT}_{X \setminus A} = -\text{DT}_A$. The distance transformation $A \mapsto \text{DT}_A$ is decreasing in the sense that $\text{DT}_A(x) \geq \text{DT}_B(x)$ for all $x \in X$ if $A \subset B$.

In the two extreme cases $A = \emptyset$ and $A = X$ we have $\text{DT}_{\emptyset} = +\infty$ and $\text{DT}_X = -\infty$. In all other cases DT_A is real-valued, $\text{DT}_A: X \rightarrow \mathbf{R}$.

Every real-valued function can be written as the difference between two nonnegative functions: $f = f^+ - f^-$, where $f^+ = \max(f, 0)$ and $f^- = \max(-f, 0)$. In particular, $\text{DT}_A = (\text{DT}_A)^+ - (\text{DT}_A)^-$. The function $(\text{DT}_A)^-$ is sometimes called the *quench function* of A .¹⁴

Proposition 5.1. *If A is a subset of a metric space X other than \emptyset and X , then $(\text{DT}_A)^+$ and $(\text{DT}_A)^-$ are Lipschitz continuous with Lipschitz constant 1 with respect to d :*

$$(5.7) \quad |(\text{DT}_A)^+(x) - (\text{DT}_A)^+(y)| \leq d(x, y), \quad x, y \in X,$$

and similarly for $(\text{DT}_A)^-$. (In particular the restrictions $\text{DT}_A|_A$ and $\text{DT}_A|_{\complement A}$ are Lipschitz continuous with Lipschitz constant 1.) As a consequence, DT_A is Lipschitz continuous with Lipschitz constant 2. If X is a vector space with distance $d(x - y) = \|x - y\|$ defined by a norm, the Lipschitz constant is 1.

¹⁴Serra (1982:377), who attributes the term to L. Calabi. However, it seems not to be widely used.

Let us say for brevity that a function is *Lip-1* if it is Lipschitz continuous with Lipschitz constant 1, i.e., when it satisfies (5.7).

Proof. The restriction to A of DT_A is the supremum of a family of Lip-1 functions $x \mapsto d(x, y)$; the restriction to $\mathfrak{C}A$ is the infimum of such a family; it is well-known that these operations preserve the Lipschitz constant. It follows that $|(\text{DT}_A)^-(a) - (\text{DT}_A)^-(b)| \leq d(a, b)$ if $a, b \in A$. If $a, b \notin A$ the function takes the value zero at both points.

Now take $a \in A$ and $b \in X \setminus A$ and define $r = -\text{DT}_A(a) \geq 0$ and $s = \text{DT}_A(b) \geq 0$. Then the strict ball $B_{<}(a, r)$ is contained in A , and the strict ball $B_{<}(b, s)$ is contained in $\mathfrak{C}A$, so that $r, s \leq d(a, b)$. The two balls are disjoint. In general this only implies that $\max(r, s) \leq d(a, b)$, but in a normed vector space case the stronger inequality $s + r \leq d(a, b)$ follows, thus that $0 \leq \text{DT}_A(b) - \text{DT}_A(a) = s + r \leq d(a, b)$; proving that the Lipschitz constant is 1 in this case.

Returning to the general case, we note that, when $a \in A$ and $b \notin A$, we have $0 = -(\text{DT}_A)^-(b) \leq \text{DT}_A(a)^- - \text{DT}_A(b)^- = r \leq d(a, b)$ and the Lipschitz continuity of $(\text{DT}_A)^-$ is established. Passing to the complement, we obtain the result for $(\text{DT}_A)^+ = (\text{DT}_{X \setminus A})^-$.

The difference $\text{DT}_A = (\text{DT}_A)^+ - (\text{DT}_A)^-$ thus has Lipschitz constant at most 2. There are easy examples which show that the constant cannot be improved.

If G is an abelian group with a translation-invariant metric, then the distance transform of any subset A can be expressed as an infimal convolution in the complement of A : $\text{DT}_A = i_A \square f$ in $G \setminus A$, where $f(x) = d(x, 0)$ is the distance from x to the neutral element 0 and i_A is the indicator function of the set. Passing to the complement we get of course $\text{DT}_A = -\text{DT}_{\mathfrak{C}A} = -(i_{\mathfrak{C}A} \square f)$ in A . We see that

$$(\text{DT}_A)^+ = \max(\text{DT}_A, 0) = i_A \square f \text{ and } (\text{DT}_A)^- = \max(-\text{DT}_A, 0) = i_{\mathfrak{C}A} \square f,$$

and, taking the difference between the two,

$$\text{DT}_A = (\text{DT}_A)^+ - (\text{DT}_A)^- = (i_A \square f) - (i_{\mathfrak{C}A} \square f).$$

Proposition 5.2. *Let G be an abelian group with a translation-invariant metric d . Then for any subsets A, B of G we have*

$$(\text{DT}_{A+B})^+ = (\text{DT}_A)^+ \square i_B = i_A \square (\text{DT}_B)^+ = (\text{DT}_A)^+ \square (\text{DT}_B)^+.$$

Proof. We know from the discussion preceding the proposition that $(\text{DT}_A)^+ = i_A \square f$, where $f(x) = d(x, 0)$ is the distance from x to the origin. Hence, using freely the associativity and commutativity of infimal convolution as well as the functional equation $f \square f = f$ (cf. Lemma 5.14 below),

$$(\text{DT}_A)^+ \square (\text{DT}_B)^+ = (i_A \square f) \square (i_B \square f) = (i_A \square i_B) \square f = (\text{DT}_{A+B})^+.$$

Also

$$(\text{DT}_A)^+ \square i_B = (i_A \square f) \square i_B = (i_A \square i_B) \square f = (\text{DT}_{A+B})^+.$$

5.2. Distance transforms and sublevel sets

The *sublevel sets* of a function $f: X \rightarrow [-\infty, +\infty]$ are the sets of the form

$$\{x \in X; f(x) < s\} \text{ or } \{x \in X; f(x) \leq s\}$$

for some element s of $[-\infty, +\infty]$. For brevity we shall denote them by $\{f < s\}$ rather than $\{x \in X; f(x) < s\}$ when no misunderstanding seems possible.

Lemma 5.3. *If X is a metric space with metric d , and DT_A is the distance transform of a subset A of X calculated with the use of d , then the closure, interior and boundary of A can all be recovered from knowledge of the sublevel sets of DT_A :*

$$\bar{A} = \{\text{DT}_A \leq 0\}, \quad A^\circ = \{\text{DT}_A < 0\}, \quad \partial A = \{\text{DT}_A = 0\}.$$

Moreover $\text{DT}_{\bar{A}} = \text{DT}_A$ in $\mathfrak{C}A$ and $\text{DT}_{A^\circ} = \text{DT}_A$ in A .

Proof. To prove the equality $\bar{A} = \{x; \text{DT}_A(x) \leq 0\}$, first note that if $x \in A$, then $x \in \bar{A}$ and $\text{DT}_A(x) \leq 0$ by definition. If on the other hand $x \notin A$, then $x \in \bar{A}$ if and only if there are points $y \in A$ such that $d(x, y)$ is arbitrarily small, which happens if and only if $\text{DT}_A(x) = 0$. This proves the first equality.

The second equality in the statement follows by passing to the complement, and the third by taking the set-theoretical difference $\bar{A} \setminus A^\circ = \partial A$.

If A is any subset of \mathbf{R}^n satisfying $B_{<}(c, r) \subset A \subset B_{\leq}(c, r)$, where we use the distance $d(x, y) = \|x - y\|$ defined by some norm $\|\cdot\|$ on \mathbf{R}^n , then $\text{DT}_A(x) = \|x - c\| - r$. This simple example shows that we cannot expect to recover A exactly from DT_A ; we have to be content with its interior and closure. However, if $X = \mathbf{Z}^n$, then the topology induced by a norm in \mathbf{R}^n is the discrete topology, so that, for any set A ,

$$\bar{A} = A^\circ = A = \{\text{DT}_A < 0\} = \{\text{DT}_A \leq 0\}.$$

Proposition 5.4. *Let G be an abelian group with a translation-invariant metric d , and let A be an arbitrary subset of G . Then for all positive numbers r and ε we have*

$$\{\text{DT}_A < r\} = A + B_{<}(0, r) \subset A + B_{\leq}(0, r) \subset \{\text{DT}_A \leq r\} \subset \{\text{DT}_A < r + \varepsilon\};$$

and

$$\{\text{DT}_A \leq -r\} = A \ominus B_{<}(0, r) \supset A \ominus B_{\leq}(0, r) \supset \{\text{DT}_A < -r\} \supset \{\text{DT}_A \leq -r - \varepsilon\}.$$

Proof. This is easy; we omit the proof. It is enough to prove the first chain of inclusions; the second follows by duality.

Give examples to show that where an inclusion sign is written, the inclusion may be strict.

The dilations $A + B_{<}(0, r)$ thus determines the strict sublevel sets of DT_A for positive values; similarly for the erosions $A \ominus B_{<}(0, r)$ and the nonstrict sublevel sets of DT_A for negative values.

Proposition 5.5. *Let G be an abelian group and $f_j: G \mapsto [-\infty, +\infty]$, $j = 1, 2$, two arbitrary functions defined on G . Define $f_3 = f_1 \square f_2$. Then for all real numbers r_1, r_2 and $r_3 = r_1 + r_2$ we have*

$$\begin{aligned} \{f_1 < r_1\} + \{f_2 < r_2\} &\subset \{f_1 < r_1\} + \{f_2 \leq r_2\} \subset \{f_3 < r_3\}; \\ \{f_1 \leq r_1\} + \{f_2 \leq r_2\} &\subset \{f_3 \leq r_3\}. \end{aligned}$$

Proof. We always have $(f_1 \square f_2)(x_1 + x_2) \leq f_1(x_1) + f_2(x_2)$, so if $f_1(x_1) < r_1$ and $f_2(x_2) \leq r_2$, then $f_3(x_1 + x_2) < r_1 + r_2 = r_3$.

In the other direction we can prove the following for distance transforms.

Proposition 5.6. *Let G be an abelian group equipped with a translation-invariant metric, and let A_j , $j = 1, 2$, be two subsets. Then their distance transforms $f_j = \text{DT}_{A_j}$ satisfy*

$$\{f_1 \square f_2 < r\} = (\{f_1 \leq 0\} + \{f_2 < r\}) \cup (\{f_1 < r\} + \{f_2 \leq 0\}) = A_1 + A_2 + B_{<}(0, r)$$

for all positive r .

What can be said when r is negative?....

Proof. That $\{f_1 \square f_2 < r\}$ contains $\{f_1 \leq 0\} \cup \{f_2 < r\}$ follows from Proposition 5.5.

To prove the inclusion in the other direction, let us assume that $(f_1 \square f_2)(x_3) < r$. Then there exists x_1, x_2 such that $x_1 + x_2 = x_3$ and $f_1(x_1) + f_2(x_2) < r$. If both $f_1(x_1)$ and $f_2(x_2)$ happen to be nonpositive, we are done. The case when $f_1(x_1)$ or $f_2(x_2)$ is positive remains to be considered. Assume first that $f_1(x_1)$ is positive. Then $x_1 \notin A_1$ and we know that, to any given positive ε , there exists a point $y_1 \in A_1$ such that $d(x_1, y_1) < f_1(x_1) + \varepsilon$. We may choose $\varepsilon = r - f_1(x_1) - f_2(x_2)$. Define $y_2 = x_1 + x_2 - y_1$. Then $f_1(y_1) \leq 0$ and the Lipschitz continuity of f_2 implies that

$$f_2(y_2) \leq f_2(x_2) + d(y_2, x_2) < f_1(x_1) + f_2(x_2) + \varepsilon = r.$$

(Note that $y_2 - x_2 = x_1 - y_1$, so that $d(y_2, x_2) = d(y_1, x_1) < f_1(x_1) + \varepsilon$.) Hence $x_3 = x_1 + x_2 = y_1 + y_2 \in \{f_1 \leq 0\} + \{f_2 < r\}$. The case $f_2(x_2) > 0$ is similar.

To prove the last equality we note that $\{f_1 \leq 0\}$ is equal to the closure $\overline{A_1}$ of A_1 and that $\{f_2 < r\} = A_2 + B_{<}(0, r)$. However $\overline{A} + B = A + B$ if B is open, so both terms in the union simplify to $A_1 + A_2 + B_{<}(0, r)$.

5.3. Distance transforms in normed vector spaces

In this section we shall calculate the distance transforms of some subsets of normed spaces. In such spaces we shall always use the metric defined by the given norm of the space, $d(x, y) = \|x - y\|$.

The space \mathbf{R}^n of all n -tuples can be normed by the l^p -norms $\|\cdot\|_p$, $1 \leq p \leq +\infty$, which are defined by

$$(5.8) \quad \|x\|_p = \left(\sum |x_j|^p \right)^{1/p}, \quad x = (x_1, \dots, x_n) \in \mathbf{R}^n.$$

When $p = +\infty$ this has to be interpreted as a limit. More explicitly one defines

$$(5.9) \quad \|x\|_\infty = \max_j |x_j|, \quad x \in \mathbf{R}^n.$$

For any normed vector space E we consider its *dual* E' , consisting of all continuous linear forms on E . On the dual we define the norm *dual* to $\|\cdot\|$ by $\|\xi\|' = \sup_{\|x\| \leq 1} |\xi(x)|$ for $\xi \in E'$. It follows that $|\xi(x)| \leq \|\xi\|' \|x\|$ for all $x \in E$ and all $\xi \in E'$.

When $E = \mathbf{R}^n$, we may identify also E' with \mathbf{R}^n , and the evaluation of ξ at the point x , i.e., the number $\xi(x)$, is then the inner product $\xi \cdot x$. The Euclidean norm $\|\cdot\|_2$ is dual to itself:

$$\|\xi\|'_2 = \sup_{\|x\|_2 \leq 1} \xi(x) = \|\xi\|_2 = \sqrt{\sum \xi_j^2}.$$

It is not difficult to prove that the norm dual to $\|\cdot\|_1$ is $\|\cdot\|_\infty$ and vice versa. More generally, one may prove that the norm dual to $\|\cdot\|_p$ is $\|\cdot\|_q$, where $q = p/(p-1)$, $1 < p < +\infty$, with a natural interpretation also when $p = 0, +\infty$. This statement follows from Hölder's inequality and its converse.

In all what follows we may take E as \mathbf{R}^n with one of these norms. However, the more general statements are really not more difficult to prove.

In particular we shall look at the case of a convex subset A of a normed vector space. We already noted that the distance transform of a ball $B = B_{\leq}(c, r)$ is $\text{DT}_B(x) = \|x - c\| - r$. Another simple convex set is a half-space, and it is easy to determine its distance transform:

Proposition 5.7. *Let Y be a closed half-space in a normed space E , defined by an inequality $\xi(x - a) \leq 0$ for some continuous linear form $\xi \in E'$, $\xi \neq 0$. Then its distance transform is $\text{DT}_Y(x) = \eta(x - a)$, where $\eta = \xi/\|\xi\|'$.*

Proof. We may as well define Y by the inequality $\eta(x - a) \leq 0$. We note that $\|\eta\|' = 1$.

By definition $\text{DT}_Y(x) = \inf_{\eta(y) \leq \eta(a)} \|x - y\|$ when $x \notin Y$. Without loss of generality we may assume that $x = 0$ and $\eta(a) = -1$; this can be achieved by a change of variable. We shall thus have to prove that

$$(5.10) \quad \inf_{\substack{y \in E \\ \eta(y) \leq -1}} \|y\| = 1.$$

This is a well-known fact. To prove it, we note that $\sup_{\|z\| \leq 1} \eta(z) = 1$ by the definition of the dual norm, and that, given any ε with $0 < \varepsilon < 1$, we can take z_ε of norm at most 1 so that $\eta(z_\varepsilon) \geq 1 - \varepsilon$. Then $y_\varepsilon = -z_\varepsilon/\eta(z_\varepsilon)$ satisfies $\eta(y_\varepsilon) = -1$ and $\|y_\varepsilon\| \leq 1/(1 - \varepsilon)$. Hence the infimum in (5.8) is at most $1/(1 - \varepsilon)$, thus at most 1. On the other hand, the infimum is at least 1, for $\|y\| \geq |\eta(y)| \geq 1$. This proves the formula for DT_Y in the complement of Y . We note that for open half-spaces the formula is the same. For $x \in Y$ we can therefore use the symmetry $\text{DT}_{\mathbf{c}Y} = -\text{DT}_Y$.

We note that if Y is a half-space which is not closed (equivalently, defined by a discontinuous linear functional), then $\text{DT}_Y = 0$ identically.

Theorem 5.8. *Let A be any convex set in a normed vector space E . Then its distance transform is*

$$\text{DT}_A = \sup_Y (\text{DT}_Y; Y \text{ is a closed half-space containing } A).$$

Proof. We first note that the theorem is trivially true if $A = \emptyset$ or $A = E$. Assume now that $A \neq \emptyset, E$.

For every set Y containing A we have $\text{DT}_Y \leq \text{DT}_A$, so the supremum in the statement of the theorem can never exceed DT_A .

If $x \notin \overline{A}$ we consider the open ball $B_{<}(x, r)$, where $r = \text{DT}_A(x) > 0$. This ball and A are disjoint. By the Hahn–Banach theorem there is a hyperplane separating the two. This means that there is a continuous linear form ξ such that $\xi(y) \leq \xi(z)$ for all $y \in A$ and all $z \in B_{<}(x, r)$. Thus the closed half-space Y defined by $\xi(y) \leq \inf_{z \in B_{<}(x, r)} \xi(z)$ contains A , and $\text{DT}_Y(x) \geq r = \text{DT}_A(x)$. This shows that the supremum in the statement of the theorem attains DT_A at the point x .

Now take $x \in A$. Then there is an open ball $B_{<}(x, r)$ with maximal radius $r = -\text{DT}_A(x) \geq 0$ contained in A . (Note that $r = 0$ is allowed this time.) Take $r' = r + \varepsilon > r$, where ε is arbitrarily small. There exists a point b in $B_{\leq}(x, r') \setminus A$. By the Hahn–Banach theorem again there is a closed hyperplane passing through b with A on one side. The corresponding half-space Y solves our problem, for $-\text{DT}_Y(x)$ cannot be larger than $d(x, b) \leq r'$, so that $\text{DT}_Y(x) \geq -r' = \text{DT}_A(x) - \varepsilon$. Since the two sides are equal in A , they are also equal in the closure of A by continuity. Since we already proved that they are equal in the complement of \overline{A} , we are done.

Corollary 5.9. *The distance transform of an arbitrary convex subset of a normed vector space is a convex function.*

Proof. The distance transform of a half-space is an affine function as we have seen. Every affine function is convex, and the supremum of any family of convex functions is convex.

Given any subset A of a vector space E we define its *supporting function* H_A by

$$(5.11) \quad H_A(\xi) = \sup_{x \in A} \xi(x), \quad \xi \in E^*,$$

where E^* is the space of all linear forms on E . It is called the *algebraic dual* of E , and contains the dual E' , maybe strictly.

Corollary 5.10. *The distance transform of a closed convex subset A of a normed space E is equal to*

$$\text{DT}_A(x) = \sup_{\|\eta\|=1} (\eta(x) - H_A(\eta)), \quad x \in E.$$

Proof. We know that the distance transform at a point x is equal to the supremum of all values $\text{DT}_Y(x)$ when Y varies in the family of all closed half-spaces containing A . However, every such half-space is contained in a minimal half-space, and the minimal half-spaces $Y = \{y; \eta(y) \leq \alpha\}$ are precisely those for which $\alpha = H_A(\eta)$, thus with distance transform $\eta(x) - H_A(\eta)$.

To any function $f: E \rightarrow [-\infty, +\infty]$ we define its *Fenchel transform* $\tilde{f}: E^* \rightarrow [-\infty, +\infty]$ on the algebraic dual E^* of E by

$$\tilde{f}(\xi) = \sup_{x \in E} (\xi(x) - f(x)), \quad \xi \in E^*.$$

The Fenchel transform generalizes the supporting function: $\tilde{i}_A = H_A$. Corollary 5.10 shows that the distance transform is the Fenchel transform of the supporting function restricted to the unit sphere. We may therefore ask what happens if we apply the transformation again.

Corollary 5.11. *If A is a nonempty closed convex subset of E , the Fenchel transform \tilde{f} of its distance transform $f = \text{DT}_A$ is*

$$\tilde{f}(\xi) = \begin{cases} H_A(\xi), & \|\xi\|' = 1; \\ +\infty, & \|\xi\|' > 1; \\ \text{the largest convex minorant of } H_A|_{S'}, & \|\xi\|' < 1. \end{cases}$$

Here S' is the unit sphere for the dual norm $\|\cdot\|'$. If A is empty, then \tilde{f} is $-\infty$ identically.

Proof. Let g denote the restriction of H_A to the unit sphere S' and h the restriction of H_A to the closed unit ball of E' . We note that, by Corollary 5.10, f is the Fenchel transform of g . The transform of f is therefore the second Fenchel transform of g : $\tilde{f} = \tilde{\tilde{g}}$. The general theory of the Fenchel transformation tells us that $\tilde{\tilde{g}}$ is the largest lower semicontinuous convex minorant of g . Since $h \leq g$ we have $\tilde{\tilde{h}} \leq \tilde{\tilde{g}}$, and we know that $\tilde{\tilde{h}} = h$ since h is convex, lower semicontinuous, and never takes the value $-\infty$; $h = \tilde{\tilde{h}} \leq \tilde{\tilde{g}} \leq g$. Therefore, since h and g agree on S' , $\tilde{\tilde{g}} = g = H_A$ on S' . Outside B' we have $g = h = +\infty$, so $\tilde{f} = \tilde{\tilde{g}} \geq \tilde{\tilde{h}} = +\infty$ there. Finally we note that, since g is plus infinity in the open unit ball, $\tilde{\tilde{g}}$ is as described in the statement of the corollary in that ball.

All this holds if A is nonempty; otherwise we easily see that $\tilde{f} = H_A = -\infty$ identically.

5.4. Chamfer distances

While the Euclidean distance is easy to visualize geometrically, it has certain drawbacks when it comes to calculations: we need to keep in memory a vector rather than a scalar at each pixel; we need more operations per pixel; and, perhaps most importantly, the Euclidean distance is more difficult to use for various morphological operations, such as skeletonizing, than for instance the city-block distance; see Borgefors (1994). For a study of the computation of the Euclidean distance transform in any dimension, see Ragnemalm (1993).

In the case of the city-block (l^1) and chessboard (l^∞) distances, one first defines the distances between neighboring pixels; we shall call them, following Starovoitov (1995:501), *prime distances*. Then the distance between any two pixels is defined by following a path and taking as the distance the minimum over all admissible paths of the sum of the prime distances. As an example, for the city-block distance the admissible paths consists of horizontal and vertical moves only, and the prime distance between two pixels which share a side is declared to be one. Thus the distance is calculated successively from neighboring pixels, which is convenient both for sequential and parallel computation. This is impossible for the Euclidean distance in spaces of dimension two or more.

It turns out that many metrics used in image analysis are conveniently defined from the prime distances by infimal convolution over all grid points. We shall now explain this, following Kiselman (1996).

The following result is well known and easy to prove.

Lemma 5.12. *Any translation-invariant distance d on an abelian group G defines a function $f(x) = d(x, 0)$ on X which is positive definite:*

$$(5.12) \quad f(x) \geq 0 \text{ with equality precisely when } x = 0;$$

and symmetric:

$$(5.13) \quad f(-x) = f(x) \text{ for all } x \in X.$$

Conversely, a function f which satisfies (5.9) and (5.10) defines a distance $d(x, y) = f(x - y)$.

Lemma 5.13. *Let d be a translation-invariant distance on an abelian group G and f a function on G related to d as in Lemma 5.12. Then d is a metric if and only if f is subadditive:*

$$(5.14) \quad f(x + y) \leq f(x) + f(y) \text{ for all } x, y \in X.$$

Proof. If d is a metric, we can write, using the triangle inequality and the translation invariance,

$$f(x + y) = d(x + y, 0) \leq d(x + y, y) + d(y, 0) = d(x, 0) + d(y, 0) = f(x) + f(y).$$

Conversely, if f is subadditive,

$$d(x, z) = f(x - z) \leq f(x - y) + f(y - z) = d(x, y) + d(y, z),$$

proving the triangle inequality.

In the definition of an infimal convolution the infimum operator acts over an infinite set of points, and therefore sometimes cannot be computed in finitely many steps. However, there are many situations where the infimum is in fact a minimum over a finite set. One such case is when f is bounded from below and g is coercive in the strong sense that all sublevel sets $\{y; g(y) \leq a\}$, $a \in \mathbf{R}$, are finite. Then in particular the sublevel set $\{y; g(y) \leq (f \square g)(x) + 1 - \inf f\}$ is finite for every x , and it is enough to search for a minimizing y in that set. Even simpler is the case when g is less than $+\infty$ in a finite set P only. Then the infimal convolution with any function f is equal to the minimum

$$(f \square g)(x) = \min_{y \in P} (f(x - y) + g(y)), \quad x \in G.$$

This is indeed the case for the distances we shall consider: here P is a small set around the origin where the prime distances are defined.

We have seen that subadditive functions are important when it comes to defining metrics (Lemma 5.2). Therefore it is of interest to know that subadditivity can be characterized using infimal convolution:

Lemma 5.14. *A function f on an abelian group is subadditive in the sense of (5.11) if and only if it satisfies the inequality $f \square f \geq f$. If $f(0) = 0$, this is equivalent to the equation $f \square f = f$.*

Proof. If f is subadditive we have $f(x - y) + f(y) \geq f(x)$, so taking the infimum over y gives $(f \square f)(x) \geq f(x)$. Conversely, $f(x) + f(y) \geq (f \square f)(x + y)$ for all x, y , so $f \square f \geq f$ implies subadditivity. Finally, we always have $(f \square f)(x) \leq f(x) + f(0)$, so if $f(0) = 0$ it follows that $f \square f \leq f$.

Infimal convolution is a commutative and associative operation on functions (see Proposition 2.1), so we can write iterated convolutions as $f \square g \square h$ without using parentheses. A k -fold convolution can be defined by

$$(5.15) \quad (f_1 \square \cdots \square f_k)(x) = \inf \sum_{i=1}^k f_i(x^i), \quad x \in G,$$

where the infimum is over all choices of elements $x^i \in G$ such that $x^1 + \cdots + x^k = x$, and with the understanding that the sum receives the value $+\infty$ as soon as one of the terms has that value, even in the presence of a value $-\infty$. In (5.12) it is natural to think of a path leading from 0 to x consisting of segments $[0, x^1], [x^1, x^1 + x^2], \dots, [x^1 + \cdots + x^{k-1}, x]$; if $G = \mathbf{Z}^2$ this path can be realized in \mathbf{R}^2 .

If A is a subset of an abelian group G , we shall write $\mathbf{N} \cdot A$ for the semigroup generated by A :

$$\mathbf{N} \cdot A = \left\{ \sum m_i a_i; m_i \in \mathbf{N}, a_i \in A \right\},$$

where all but finitely many of the m_i are zero. Similarly, we shall write $\mathbf{Z} \cdot A$ for the group generated by A :

$$\mathbf{Z} \cdot A = \left\{ \sum m_i a_i; m_i \in \mathbf{Z}, a_i \in A \right\}.$$

If A is symmetric, $A = -A$, then of course $\mathbf{Z} \cdot A = \mathbf{N} \cdot A$.

It seems plausible that if a repeated convolution $F \square F \square \cdots \square F$ has a limit f as the number of terms tends to infinity, then this limit will satisfy the equation $f \square f = f$. Indeed this is the case under very general hypotheses:

Theorem 5.15. *Let $F: G \rightarrow [0, +\infty]$ be a function on an abelian group G satisfying $F(0) = 0$. Define a sequence of functions $(F_j)_{j=1}^\infty$ by putting $F_1 = F$, $F_j = F_{j-1} \square F$, $j = 2, 3, \dots$, in other words, F_j is the infimal convolution of j terms all equal to F . Then the sequence $(F_j)_j$ is decreasing and its limit $\lim F_j = f \geq 0$ is subadditive. Moreover $\text{dom } f = \mathbf{N} \cdot \text{dom } F$, i.e., f is finite precisely in the semigroup generated by $\text{dom } F$.*

Remark. It is easy to prove that f is the largest subadditive minorant of F .

Proof. That the sequence is decreasing is obvious if we take $y = 0$ in the definition of F_{j+1} :

$$F_{j+1}(x) = \inf_y (F_j(x - y) + F(y)) \leq F_j(x) + F(0) = F_j(x).$$

Next we shall prove that $f(x+y) \leq f(x) + f(y)$. If one of $f(x), f(y)$ is equal to $+\infty$ there is nothing to prove, so let x, y be given with $f(x), f(y) < +\infty$ and fix a positive number ε . Then there exist numbers j, k such that $F_j(x) \leq f(x) + \varepsilon$ and $F_k(y) \leq f(y) + \varepsilon$. By associativity $F_{j+k} = F_j \square F_k$, so we get

$$f(x+y) \leq F_{j+k}(x+y) \leq F_j(x) + F_k(y) \leq f(x) + f(y) + 2\varepsilon.$$

Since ε is arbitrary, the inequality $f(x+y) \leq f(x) + f(y)$ follows. Finally, the statement about $\text{dom } g$ is an easy consequence of (2.7).

Theorem 5.16. *With F as in Theorem 5.15, assume in addition that there is a translation-invariant metric d_1 on G such that $F(x) \geq d_1(x, 0)$ for all $x \in G$. Then the limit f of the sequence F_j also satisfies this inequality, $f(x) \geq d_1(x, 0)$, so that it is positive definite. If F is symmetric, f is also symmetric and defines a metric $d(x, y) = f(x - y) \geq d_1(x, y)$ on the subgroup $\mathbf{Z} \cdot P = \mathbf{N} \cdot P$ of G generated by $P = \text{dom } F$.*

Proof. Define $H(x) = d_1(x, 0)$ and let H_j be the infimal convolution of j terms equal to H . From Lemmas 5.13 and 5.14 it follows that $H \square H = H$ and so all H_j are equal to H . Therefore $F \geq H$ implies $F_j \geq H$ and also the limit f must satisfy $f \geq H$. This proves the theorem.

When applying this theorem we could for instance let d_1 be εd_0 , where ε is a small positive number and d_0 is the discrete metric defined by (5.4). In \mathbf{Z}^n we can also use $d_1(x, y) = \varepsilon \|x - y\|$ for any norm on \mathbf{R}^n .

Corollary 5.17. *Let P be a finite set in an abelian group G containing the origin, and let F be a function on G with $F(0) = 0$, taking the value $+\infty$ outside P and finite positive values at all points in $P \setminus \{0\}$. Then $f = \lim F_j$ is a positive definite subadditive function. If P is symmetric and $F(-x) = F(x)$, then f defines a metric on the subgroup $\mathbf{Z} \cdot P = \mathbf{N} \cdot P$ of G generated by P .*

Proof. Since P is finite, there is a positive number ε such that $F(x) \geq \varepsilon$ for all $x \in P$ except $x = 0$. Thus $F(x) \geq \varepsilon d_0(x, 0)$, where d_0 is the discrete metric defined by (5.4). We can now apply the theorem.

It is by no means necessary that f be positively homogeneous. In fact, we can let $P = \{0, \pm 1, \pm 2\} \subset \mathbf{Z}$ and define $F(\pm 1) = a$, $F(\pm 2) = b$, where a and b are arbitrary positive numbers. If $b \geq 2a$, then $f(x) = a|x|$ for all $x \in \mathbf{Z}$, but if $b < 2a$, then $f(x) = \frac{1}{2}b|x|$ when x is even, $x = 2k$, $k \in \mathbf{Z}$, whereas $f(x) = b|k| + a > \frac{1}{2}b|x|$ when x is odd, $x = \pm(2k + 1)$, $k \in \mathbf{N}$. Nevertheless f is subadditive.

A more interesting example is perhaps this in two dimensions. Let $P = \{x \in \mathbf{Z}^2; |x_j| \leq 1\}$, and define the prime distances as $F(\pm 1, 0) = F(0, \pm 1) = a > 0$, $F(\pm 1, \pm 1) = b > 0$. Then if $b \geq a$ we get $f(x_1, 0) = a|x_1|$. But if $b < a$, then $f(2, 0) = 2b < 2a$, so that $f(2, 0) < 2f(1, 0) = 2a$. In fact, by the definition of infimal convolution, $f(2, 0) \leq F_2(2, 0) \leq F(1, 1) + F(1, -1) = b + b$. On the other hand, it is not difficult to see that for any k , $F_k(2, 0) \geq 2b$, so that actually $f(2, 0) = 2b$. This is because if we take $k \geq 2$ nonzero steps to go from the origin to $(2, 0)$, the distance assigned to the path is at least $F(x^1) + \dots + F(x^k) \geq kb$.

Several metrics on \mathbf{Z}^2 have been studied. When presenting the generating function F defining the prime distances it shall be understood in the sequel that F is invariant under permutation and reflection of the coordinates. Therefore it is enough to define $F(x)$ for $0 \leq x_2 \leq x_1$. Also it is understood that $F(0) = 0$ in all cases, and that $F(x) = +\infty$ when not mentioned. Consider first $P = \{x \in \mathbf{Z}^2; \sum |x_j| \leq 1\}$ and $F(1, 0) = 1$. Then the corresponding metric is the city-block metric, introduced and studied by Rosenfeld & Pfaltz (1966). If instead we let $P = \{x \in \mathbf{Z}^2; |x_j| \leq 1\}$ and $F(1, 0) = F(1, 1) = 1$, then the metric is the chessboard metric, introduced by Rosenfeld & Pfaltz (1968). Some other metrics that have been studied are modifications of this; to define them, put $F(1, 0) = a$ and $F(1, 1) = b$. Then the choices $(a, b) = (1, \sqrt{2})$ (Montanari (1968)); $(a, b) = (2, 3)$ (Hilditch & Rutovitz (1969)); and $(a, b) = (3, 4)$ (Borgefors (1984)) have all been studied. Next we can increase the size of the neighborhood where prime distances are defined to include the knight's move $(2, 1)$ as an element of P . The distance defined by this move only has been studied by Das & Chatterji (1988). It seems more natural, however, to allow also $(1, 0)$ and $(1, 1)$ in P . Then a very good choice under certain criteria is $F(1, 0) = 5$, $F(1, 1) = 7$, and $F(2, 1) = 11$ (the 5-7-11 weighted distance). This distance was proposed and studied by Borgefors (1986).

We always have $f \leq F$, and it may happen that $f(x) < F(x)$ for some pixel $x \in P$. Let for instance $F(1, 0) = a$, $F(2, 1) = c$, and extend F by reflection and permutation of the coordinates. Then

$$f(1, 0) \leq F_3(1, 0) \leq F(2, 1) + F(1, -2) + F(-2, 1) = 3c,$$

so if $3c < a$ we get $f(1, 0) \leq 3c < a = F(1, 0)$. This is undesirable, because we expect the prime distance originally defined between the origin and $(1, 0) \in P$ to survive and to be equal to the distance defined by the minimum over all paths. It is therefore natural to require that $f = F$ everywhere in P .

Definition. 5.18. *Let us say that a metric $d(x, y) = f(x - y)$ is a chamfer distance, or finitely generated if it is constructed as in Corollary 5.17.*

It is easy to prove that the Euclidean metric $d(x, y) = \sqrt{\sum (x_j - y_j)^2}$ on \mathbf{Z}^n is a chamfer distance if and only if $n \leq 1$.

Borgefors (1984:324, 1986:345) calls $F(x - y)$ the *local distances*; Verwer (1991:672) adopted this term. Starovoitov (1995:501) calls F the *prime distance function*. The resulting metric, $(x, y) \mapsto f(x - y)$, was called a *quasi-Euclidean* distance by Montanari (1968) and a *chamfer distance* by Borgefors (1984:326). The term has then been used, e.g., by Verwer (1991:672) and Marchand-Maillet & Sharaiha (2000:21). However, Borgefors herself now prefers the term *weighted distances*, and refers to them as *constructed by chamfering*;¹⁵ the latter term is derived from the method of calculating distance transforms by passing a mask twice over the image. Nevertheless, the term *chamfer distance* for the metrics constructed in Corollary 5.17 has won acceptance.

Since $f(x)$ is the limit of an infinite sequence $F_j(x)$, it is reassuring to know that this sequence is in fact stationary in the cases of interest here. It is easy to explicitly give an index j such that $F_j(x)$ is equal to the limit $f(x)$:

¹⁵Gunilla Borgefors, personal communication 2002-03-22.

Proposition 5.19. *Let F be as in Corollary 5.17. Then the sequence (F_j) is pointwise stationary, i.e., for every $x \in G$ there is an index p_x such that $F_j(x) = f(x)$ for all $j \geq p_x$.*

Proof. If $f(x) = +\infty$, then of course $F_j(x) = +\infty$ for all j . Consider a point such that $f(x) < \infty$. Then there is an index m such that $F_m(x) < +\infty$. For every j there are elements y^1, y^2, \dots, y^j in P such that $x = y^1 + \dots + y^j$ and $F_j(x) = F(y^1) + \dots + F(y^j)$ (cf. (5.12)). Now by hypothesis $F(y) \geq \varepsilon > 0$ for all $y \in P$ except $y = 0$, so that $F_j(x) \geq p_x \varepsilon$, where p_x is the number of indices i such that $y^i \neq 0$. If $j \geq m$ we conclude that $F_m(x) \geq F_j(x) \geq p_x \varepsilon$, so that $p_x \leq F_m(x)/\varepsilon$. Therefore the number p_x of nonzero terms in a representation of $F_j(x)$ is never larger than $F_m(x)/\varepsilon$. Now $F_j(x) = F_{p_x}(x)$, so the fact that p_x is bounded means that the sequence is stationary.

Summing up, we note that the positive part of a distance transform is a limit $(DT_A)^+ = \lim(i_A \square F \square F \square \dots \square F)$, where the number of terms tends to infinity. This formula can be used in actual calculations: starting from $g_0 = i_A$ one calculates $g_j = g_{j-1} \square F$ and stops when the criterion of Proposition 5.18 is satisfied. (It is not true that $g_j(x) = g_{j-1}(x)$ implies that $g_k(x) = g_j(x)$ for all $k \geq j$, so $g_j(x) = g_{j-1}(x)$ at a particular point x is not a sufficient criterion. Example)

5.5. Comparing distances

The l^1 (city-block) and l^∞ (chessboard) metrics are translation invariant but not rotation invariant. (In the plane a rotation can distort distances by a factor of up to $\sqrt{2}$; in higher dimensions more.) The Euclidean metric is rotation invariant, and it is desirable to construct a chamfer distance in \mathbf{Z}^n which is reasonably close to being rotation invariant. There are many studies on the problem of defining an optimal distance in a given family of finitely generated distances. Of course the property of being optimal depends on the criteria employed; as always, beauty is in the eye of the beholder. A basic problem is how to measure deviation: we may ask how far the quotient of two quantities is from 1, alternatively how far their difference is from 0. In this section we shall look briefly into this problem and describe four methods of comparing two nonnegative functions; two of these have been studied in the literature.

It is natural to measure the deviation of a function $f: X \rightarrow [0, +\infty[$ from a given nonnegative function g defined on the same set by the smallest constant $C \in [0, +\infty[$ such the inequalities $f(x) \leq Cg(x)$ and $g(x) \leq Cf(x)$ hold for all $x \in X$. We introduce a notation for this constant,

$$(5.16) \quad \Lambda(f, g) = \max \left(\sup_{x \in X} \frac{f(x)}{g(x)}, \sup_{x \in X} \frac{g(x)}{f(x)} \right),$$

where the supremum is taken over all points in the common domain of definition, and where we count $0/0$ as 0 and $t/0$ as $+\infty$ if $t > 0$ (this is to allow for zeros; $\Lambda(f, g)$ is finite only if the two functions have the same zero set). It is noteworthy that $\log \Lambda(f, g) = \|\log f - \log g\|_\infty$ is a distance on a suitable space of functions.

If f satisfies an inequality $C_1 \leq f(x)/g(x) \leq C_2$, then a slightly modified function, viz. $f_1 = f/\sqrt{C_1 C_2}$, satisfies $\Lambda(f_1, g) \leq \sqrt{C_2/C_1}$.

Another, closely related measure of the deviation was studied by Verwer (1991).

He used the functional

$$(5.17) \quad \Lambda'(f, g) = \sup_{x \in X} \left| \frac{f(x)}{g(x)} - 1 \right|.$$

However, one might just as well consider $\Lambda'(g, f)$. It is readily seen that $\Lambda'(f, g) = \Lambda(f, g) - 1$ when $f \geq g$, and that $\Lambda'(f, g) = 1 - 1/\Lambda(f, g)$ when $f \leq g$. In general, $\Lambda'(f, g)$, $\Lambda'(g, f)$ as well as $\log \Lambda(f, g)$ lie between two limits,

$$1 - \frac{1}{\Lambda(f, g)} \leq \Lambda'(f, g), \Lambda'(g, f), \log \Lambda(f, g) \leq \Lambda(f, g) - 1,$$

where we have put in also the well-known inequality $1 - 1/t \leq \log t \leq t - 1$. In particular, $|\Lambda'(f, g) - \log \Lambda(f, g)| \leq (\Lambda(f, g) - 1)^2/\Lambda(f, g)$; the same estimate holds of course for $\Lambda'(g, f)$. We may also note that, although Λ' is not symmetric, it is approximately symmetric when f and g are close, and there is an estimate

$$\frac{\Lambda'(f, g)}{\Lambda(f, g)} \leq \Lambda'(g, f) \leq \Lambda(f, g)\Lambda'(f, g),$$

which, by the way, may be written as

$$\Lambda(\Lambda'(f, g), \Lambda'(g, f)) \leq \Lambda(f, g).$$

When f and g are reasonably close, $\Lambda'(f, g) \approx \Lambda'(g, f) \approx \log \Lambda(f, g)$. For many purposes either one may be used. Note, however, that $\Lambda(f, g)$ has better functional properties than $\Lambda'(f, g)$. In particular, as already noted, $\log \Lambda(f, g)$ is a metric, whereas $\Lambda'(f, g)$ does not satisfy the triangle inequality and is not even symmetric.

We note that for every pair (f, g) of functions there are constants c_0 , c_1 , and c_2 such that, respectively, $\Lambda(c_0f, g)$, $\Lambda'(c_1f, g)$ and $\Lambda'(g, c_2f)$ are minimal. It is easy to see that c_0 is the geometric mean of c_1 and c_2 .

In particular we shall compare the chamfer distances with the Euclidean norm $g = \|\cdot\|_2$. In the exercises below we let the reader determine or estimate the deviation from the Euclidean norm of some well-known finitely generated metrics. If the prime vectors are $(\pm 1, 0)$, $(0, \pm 1)$ and $(\pm 1, \pm 1)$ with prime distances a and b respectively, we note that the optimal prime distances for both Λ and Λ' are related by $b = a\sqrt{2}$, but that the actual values are slightly different. For $\Lambda(f, \|\cdot\|_2)$, the optimal choice is

$$(5.18) \quad a = a_0 = \sqrt[4]{\frac{2 + \sqrt{2}}{4}} \approx 0.961186523, \quad b = b_0 = a_0\sqrt{2} = \sqrt[4]{2 + \sqrt{2}} \approx 1.359323017,$$

whereas Verwer (1991:676) found the optimal choice for $\Lambda'(f, \|\cdot\|_2)$ to be

$$a = a_1 \approx 0.9604 \text{ and } b = b_1 = a_1\sqrt{2} \approx 1.3583.$$

The exact values are

$$(5.19) \quad a_1 = \frac{1}{\frac{1}{2} + \sqrt{1 - 1/\sqrt{2}}} \text{ and } b_1 = a_1\sqrt{2}.$$

One can calculate also the optimal choice for $\Lambda'(\|\cdot\|_2, f)$, which is

$$(5.20) \quad a = a_2 = \frac{1}{2} + \frac{1}{4}\sqrt{2 + \sqrt{2}} \approx 0.961939766 \text{ and } b = b_2 = a_2\sqrt{2} \approx 1.3603882.$$

In this case the optimality has a clear geometrical meaning: the vertices of the octagon protrude as much outside the disk as the midpoints of the edges go into the disk. As we can expect from an earlier remark, $a_0 = \sqrt{a_1 a_2}$.

In her pioneering work on finitely generated distances (1984, 1986), Gunilla Borgfors used instead the functional

$$(5.21) \quad \Lambda'_{2,\infty}(f) = \sup_{x \in \mathbf{Z}^2} \left| \frac{f(x)}{\|x\|_2} - 1 \right| \frac{\|x\|_2}{\|x\|_\infty} = \sup_{x \in \mathbf{Z}^2} \frac{|f(x) - \|x\|_2|}{\|x\|_\infty},$$

and she determined the optimal distances for this measure as follows:

$$\Lambda'_{2,\infty}(f) = \frac{1}{2}\sqrt{2\sqrt{2} - 2} - \frac{1}{2} \approx 0.04491,$$

attained for

$$(5.22) \quad a = a_3 = \frac{1}{2}\sqrt{2\sqrt{2} - 2} + \frac{1}{2} \approx 0.95509 \text{ and } b = b_3 = \sqrt{2} - \frac{1}{2} + \frac{1}{2}\sqrt{2\sqrt{2} - 2} \approx 1.36930;$$

note that $b_3 > a_3\sqrt{2}$ here (1986:351). She also determined the optimal values for the functional $\Lambda'_{2,\infty}$ when a is restricted to be 1 and b is free to vary; in this case $\Lambda'_{2,\infty}(f) = 1/\sqrt{2} - \sqrt{\sqrt{2} - 1} \approx 0.06$ and is attained for

$$(5.23) \quad a = a_4 = 1 \text{ and } b = b_4 = 1/\sqrt{2} - \sqrt{\sqrt{2} - 1} \approx 1.351;$$

here $b_4 < a_4\sqrt{2}$ (1984:327).

5.6. The calculus of balls

In any metric space the inequality $d(c, c') + r \leq r'$ implies that the open ball $B_{<}(c, r)$ is contained in $B_{<}(c', r')$ and also that $B_{\leq}(c, r) \subset B_{\leq}(c', r')$. In a normed vector space of dimension at least one and equipped with the distance $d(x, y) = \|x - y\|$ defined by the norm, the converse implications hold, provided $r > 0$ ($r \geq 0$ for closed balls). In particular, two open balls $B_{<}(c, r)$ and $B_{<}(c', r')$ with $r > 0$ are equal if and only if $c = c'$ and $r = r'$ ($r \geq 0$ suffices for closed balls). But in a general metric space very little can be said: from $B_{\leq}(c, r) \subset B_{\leq}(c', r')$ we can only deduce that $d(c, c') \leq r'$, assuming that $r \geq 0$.

We note, however, that if, in an abelian group with a translation-invariant metric, $B_{\leq}(c, r)$ is contained in $B_{\leq}(c', r')$, $r' \leq r$, and $r \geq 0$, then $c - c'$ belongs to a bounded subgroup. In the most common applications, like normed spaces and groups with a chamfer distance, the only bounded subgroup is $\{0\}$, thus $c = c'$.

The mappings

$$X \times \mathbf{R} \ni (c, r) \mapsto B_{\leq}(c, r), B_{<}(c, r) \in \mathcal{P}(X)$$

are in general far from injective and induce complicated equivalence relations in $X \times \mathbf{R}$. Similarly, the inclusion relations $B_{\leq}(c, r) \subset B_{\leq}(c', r')$, $B_{<}(c, r) \subset B_{<}(c', r')$ induce

preorders in $X \times \mathbf{R}$ (see (3.1) and (3.2)). The purpose of this section is to study these relations.

In any metric space X with metric d we fix a point c in X and define two functions which are of interest when defining balls:

$$\rho(r) = \sup_{x \in X} (d(x, c); d(x, c) \leq r) \text{ and } \sigma(r) = \inf_{x \in X} (d(x, c); d(x, c) \geq r), \quad r \in [-\infty, +\infty].$$

The functions depend of course in general on the choice of c , but in most applications we shall let X be an abelian group and assume that d is translation invariant. Then the choice of c is immaterial.

The functions

$$\rho_-(r) = \sup_{s < r} \rho(s) = \sup_{x \in X} (d(x, c); d(x, c) < r), \quad r \in [-\infty, +\infty],$$

and

$$\sigma_+(r) = \inf_{s > r} \sigma(s) = \inf_{x \in X} (d(x, c); d(x, c) > r), \quad r \in [-\infty, +\infty],$$

are also of interest. It is clear that

$$\rho_- \leq \rho \leq \text{Id}_{[-\infty, +\infty]} \leq \sigma \leq \sigma_+.$$

On the negative half-axis we have $\rho(r) = -\infty$ for $-\infty \leq r < 0$ and $\sigma(r) = 0$ for $-\infty \leq r \leq 0$.

We have $B_{\leq}(c, r) = B_{\leq}(c, \rho(r))$ and $B_{<}(c, r) = B_{<}(c, \sigma(r))$. The interval $]\rho(r), \sigma(r)[$ contains no distances $d(c, x)$ —a sphere with center at c and radius in that interval is empty. The functions ρ and σ solve the uniqueness problem for balls with a common center: $B_{\leq}(c, r) = B_{\leq}(c, r')$ if and only if $\rho(r) = \rho(r')$, and $B_{<}(c, r) = B_{<}(c, r')$ if and only if $\sigma(r) = \sigma(r')$.

The mappings $\rho, \sigma: [-\infty, +\infty] \rightarrow [-\infty, +\infty]$ are idempotent; in fact ρ is an opening and σ is a closing (Definitions 3.1 and 3.3). By way of contrast, ρ_- and σ_+ are in general not idempotent. However, $\rho \circ \rho_- = \rho_-$.

In a metric space with a translation-invariant distance, the value of the distance transform of a set A at a point $x \in A$, $\text{DT}_A(x) = -r$, satisfies $r = \sigma(r)$. Similarly for $x \notin A$, we have $\text{DT}_A(x) = r$ with $r = \rho(r)$. Thus the values of $|\text{DT}_A|$ are always contained in the invariant set Inv_σ of σ .

In a normed vector space of dimension at least one, we always have $\rho(r) = \sigma(r) = r$ for positive r , so there is no need to introduce them in the study of such spaces. In \mathbf{Z}^n equipped with the l^∞ or l^1 metrics we have $\rho(r) = \lfloor r \rfloor$ and $\sigma(r) = \lceil r \rceil$ for positive r .

While the functions ρ and σ solve the problem of comparing two balls with a common center, more hypotheses will be needed if we are to compare successfully balls with different centers. This is why we are led to an analysis of the triangle inequality.

We shall say that a translation-invariant metric $d(x, y) = f(x - y)$ defined on an abelian group is *upper regular for the triangle inequality* if, given any x and y , there is a point \tilde{y} such that

$$(5.24) \quad f(\tilde{y}) = f(y) \text{ and } f(x + \tilde{y}) = f(x) + f(\tilde{y}).$$

We shall say that the distance is *lower regular for the triangle inequality* if given any x and y such that $f(x) \geq f(y)$, there is a point \tilde{x} such that

$$(5.25) \quad f(\tilde{x}) = f(x) \text{ and } f(\tilde{x}) = f(\tilde{x} - y) + f(y).$$

Upper regularity means that in the triangle inequality $f(x + y) \leq f(x) + f(y)$ we can always find \tilde{y} at the same distance from the origin as y and turning the inequality into an equality (by raising the left-hand side); lower regularity means that we can lower the right-hand side in the triangle inequality $f(x) \leq f(x - y) + f(y)$ without changing the left-hand side by replacing x by \tilde{x} —but of course only if the left-hand side is at least as large as $f(y)$.

It is clear that normed spaces are both upper and lower regular for the triangle inequality; we may choose \tilde{y} as a suitable multiple of x if $x \neq 0$ and \tilde{x} as a suitable multiple of y if $y \neq 0$. It is also easy to see that in \mathbf{Z}^n the l^∞ and l^1 distances are upper and lower regular for the triangle inequality.

The chamfer distance in \mathbf{Z}^2 with $P = \{0, (\pm 1, 0), (0, \pm 1), (\pm 1, \pm 1)\}$ is upper regular provided $a = F(1, 0) \leq b = F(1, 1)$. If $b \geq 2a$, then $f(x) = \|x\|_1$; if $b = a$, then $f(x) = \|x\|_\infty$. In the remaining cases, when $a < b < 2a$, the balls are octagonal, and we may argue as follows. By symmetry and reflection of the coordinates it is enough to consider a point $x \in \mathbf{Z}^2$ with $0 \leq x_2 \leq x_1$. Then x can be written as $m(1, 0) + n(1, 1)$ for some uniquely determined integers $m, n \geq 0$, and $f(x) = ma + nb$. If y is in the same sector, thus satisfying $0 \leq y_2 \leq y_1$, then $f(x + y) = f(x) + f(y)$, so we may take $\tilde{y} = y$. If y is in one of the other seven sectors we take $\tilde{y} = (|y_1|, |y_2|)$ or $(|y_2|, |y_1|)$ so that \tilde{y} is in the same sector as x and we know that we have equality in the triangle inequality (5.24). However, this chamfer distances is not lower regular if $a < b < 2a$ as can be shown by simple examples. It seems therefore that lower regularity is too stringent a criterion.

Theorem 5.20. *Let G be an abelian group and $d(x, y) = f(x - y)$ a translation-invariant metric in G . Assume that d is upper regular for the triangle inequality. Then for all $r, r' \in \mathbf{R}$, $B_{\leq}(c, r)$ is contained in $B_{\leq}(c', r')$ if and only if $d(c, c') + \rho(r) \leq \rho(r')$.*

Proof. By the definition of ρ , $B_{\leq}(c, r) = B_{\leq}(c, \rho(r))$. The inequality $d(c, c') + \rho(r) \leq \rho(r')$ implies that $B_{\leq}(c, \rho(r)) \subset B_{\leq}(c, \rho(r'))$ and hence that $B_{\leq}(c, r) \subset B_{\leq}(c', r')$.

For the other implication we argue as follows, assuming that $B_{\leq}(0, r)$ is contained in $B_{\leq}(c', r')$ (we may take $c = 0$ to simplify notation). Take a point x in the first ball. Then $f(x) \leq \rho(r)$. It follows that $f(x - c') \leq \rho(r')$. In the triangle inequality $f(x - c') \leq f(x) + f(c')$ we can find an element \tilde{x} with the same distance to the origin as x and turning the inequality into an equality: $f(\tilde{x} - c') = f(\tilde{x}) + f(c')$. Now $f(\tilde{x}) = f(x)$ can take values as close to $\rho(r)$ as we like, say larger than $\rho(r) - \varepsilon$ for a given positive ε . We get $\rho(r') \geq f(\tilde{x} - c') = f(\tilde{x}) + f(c') \geq \rho(r) - \varepsilon + d(c, c')$, proving the inequality. Note that $f(\tilde{x} - c') \leq \rho(r')$ follows from the inequality $f(\tilde{x}) = f(x) \leq \rho(r)$ and the inclusion $B_{\leq}(0, r) \subset B_{\leq}(c', r')$.

Theorem 5.21. *Let G be an abelian group and $d(x, y) = f(x - y)$ a translation-invariant metric in G . Assume that d is lower regular for the triangle inequality. Then for all $r, r' \in \mathbf{R}$ with $r > 0$, $B_{<}(c, r)$ is contained in $B_{<}(c', r')$ if and only if $d(c, c') + \sigma(r) \leq \sigma(r')$.*

Proof. By the definition of σ , $B_{<}(c, r) = B_{<}(c, \sigma(r))$. The inequality $d(c, c') + \sigma(r) \leq \sigma(r')$ implies that $B_{<}(c, \sigma(r)) \subset B_{<}(c', \sigma(r'))$, and hence that $B_{<}(c, r) \subset B_{<}(c', r')$.

To prove the other implication we assume that $B_{<}(c, r) \subset B_{<}(0, r')$ (we may take $c' = 0$). This means that $f(x) \geq \sigma(r')$ implies $f(x - c) \geq \sigma(r)$. For any given $\varepsilon > 0$ we can find x such that $\sigma(r') \leq f(x) \leq \sigma(r') + \varepsilon$. Then $f(x - c) \geq \sigma(r)$ and we obtain the rather useless string of inequalities

$$\sigma(r) \leq f(x - c) \geq f(x) - f(c) \leq \sigma(r') + \varepsilon - f(c).$$

However, since $f(c) < r' \leq \sigma(r') \leq f(x)$ (at this point we need to know that $c \in B_{<}(c, r)$, thus that r is positive), we can by hypothesis find \tilde{x} such that $f(\tilde{x}) = f(x)$ (thus $f(\tilde{x} - c) \geq \sigma(r)$) and turning the triangle inequality \geq in the last formula into an equality:

$$\sigma(r) \leq f(\tilde{x} - c) = f(\tilde{x}) - f(c) \leq \sigma(r') + \varepsilon - f(c).$$

If we now let ε tend to zero we obtain $\sigma(r) \leq \sigma(r') - f(c)$, which concludes the proof.

Since chamfer distances are in general not lower regular for the triangle inequality it is desirable to impose only upper regularity. Now upper regularity was used for the inclusion of closed balls, lower regularity for the inclusion of open balls. But, luckily, for chamfer distances all open balls are closed! We shall now study such spaces.

We consider metric spaces X such that the set of all distances $\{d(x, c); x \in X\}$ from a fixed point $c \in X$ is a discrete subset of \mathbf{R} . This is equivalent to the statement that the set $\{d(x, c); x \in X, d(x, c) \leq t\}$ is finite for every real number t . The main examples are of course the chamfer distances in any abelian group.

If the set $\{d(x, c); x \in X\}$ of all distances from a fixed point c is discrete, the set is denumerable and either finite or equal to a set $\{r_j; j \in \mathbf{N}\}$. In the latter case we shall choose the indices so that $r_0 = 0$ and $r_j < r_{j+1}$ for all $j \in \mathbf{N}$. A sphere $S(c, r)$ is nonempty if and only if $r = r_j$ for some j . The functions ρ and σ defined above satisfy

$$\rho(r) = r_j, \text{ when } r \in [r_j, r_{j+1}[, \quad \rho_-(r) = r_j, \text{ when } r \in]r_j, r_{j+1}],$$

and

$$\sigma(r) = r_j, \text{ when } r \in]r_{j-1}, r_j], \quad \sigma_+(r) = r_j, \text{ when } r \in [r_{j-1}, r_j[.$$

Proposition 5.22. *Let X be a metric space where the set of all distances from a fixed point c is discrete. Then all closed balls with center at c are open and all open balls with center at c are closed. More exactly,*

$$B_{\leq}(c, r) = B_{<}(c, \sigma_+(r)) \text{ and } B_{<}(c, r) = B_{\leq}(c, \rho_-(r)).$$

Proof. We consider the case when the set of distances is infinite—easy modifications will take care of the case of finitely many distances. It is convenient in the proof to define $r_{-1} = -\infty$. Given any $r \in \mathbf{R}$ we choose $j \in \mathbf{N}$ so that $r \in [r_{j-1}, r_j[$. Then $B_{\leq}(c, r) = B_{<}(c, s)$ for all $s \in]r_{j-1}, r_j]$. In particular $B_{\leq}(c, r) = B_{<}(c, r_j) = B_{<}(c, \sigma_+(r))$.

Similarly, given any $r \in \mathbf{R}$ we choose $j \in \mathbf{N}$ such that $r \in]r_{j-1}, r_j]$; then $B_{<}(0, r) = B_{\leq}(0, s)$ for all $s \in [r_{j-1}, r_j[$. In particular, $B_{<}(c, r) = B_{\leq}(c, r_{j-1}) = B_{\leq}(c, \rho_-(r))$ for all $r \in \mathbf{R}$.

Theorem 5.23. *Let G be an abelian group and $d(x, y) = f(x - y)$ a translation-invariant metric in G . Assume that d is upper regular for the triangle inequality and that the set of all distances is discrete. Then for all $r, r' \in \mathbf{R}$, $B_{<}(c, r)$ is contained in $B_{<}(c', r')$ if and only if $d(c, c') + \rho_-(r) \leq \rho_-(r')$.*

Proof. We know that the open ball $B_{<}(c, r)$ is equal to the closed ball $B_{\leq}(c, \rho_-(r))$. Theorem 5.19 shows that the inclusion $B_{\leq}(c, \rho_-(r)) \subset B_{\leq}(c', \rho_-(r'))$ is equivalent to the inequality $d(c, c') + \rho(\rho_-(r)) \leq \rho(\rho_-(r'))$. However, we always have $\rho(\rho_-(r)) = \rho_-(r)$. This completes the proof.

A final remark on chamfer distances. The elements of the finite set P will be denoted by $p_0 = 0$, $p_1 = -p_{-1}, \dots$, $p_m = -p_{-m}$, and the values of the prime distance function F at $\pm p_j = p_{\pm j}$, $j = 1, \dots, m$, by $a_j > 0$. It is of interest to consider the functions $\varphi, \psi: [-\infty, +\infty] \rightarrow [-\infty, +\infty]$ defined by

$$\varphi(r) = \sup_{n_j \in \mathbf{N}} \left(\sum n_j a_j; \sum n_j a_j < r \right) \text{ and } \psi(r) = \sup_{n_j \in \mathbf{N}} \left(\sum n_j a_j; \sum n_j a_j \leq r \right).$$

The supremum in the definition of ρ is over all distances from the origin that can occur; these are all of the form of a finite sum $\sum_1^m n_j a_j$ with $n_j \in \mathbf{N}$, which shows that $\rho(r) \leq \psi(r) \leq r$. Also $\rho_-(r) \leq \varphi(r) < r$. However, it may happen that $\rho(r) < \psi(r)$ for some values of r , because a certain sum $\sum n_j a_j$, although it is allowed in the definition of an infimal convolution $F \square \dots \square F$, never occurs as an infimum. The property $\psi = \rho$ is an important regularity property....??

We note that we have $\psi(r) = \rho(r)$ for \mathbf{Z}^n with the l^∞ or l^1 metrics. Indeed, here all a_j are equal to 1 and all integers can occur as distances $\sum n_j a_j = \sum n_j$. The same is true for the example in \mathbf{Z}^2 we considered with P consisting of nine points and $a = F(1, 0)$ and $b = F(1, 1)$ provided $a \leq b \leq 2a$; these inequalities guarantee that all numbers of the form $ma + nb$, $m, n \in \mathbf{N}$, do occur as distances....?? [Compare with semiregularity in the sense of Borgefors (1994, Definition 2)....??]

5.7. Exercises

5.1. Let A be the (filled) square $\{x \in \mathbf{R}^2; |x_1|, |x_2| \leq 1\}$ and B a Euclidean closed disk of radius r , viz. $B = \{x \in \mathbf{R}^2; x_1^2 + x_2^2 \leq r^2\}$. Describe for all positive r the seven sets $A + B$, $A \ominus B$, $B \ominus A$, A^B , B^A , A_B , B_A . Calculate the distance transforms DT_A and DT_B when the distance is given by $d(x, y) = \|x - y\|_p$, $p = 1, 2, \infty$.

5.2. Given a subset A of a normed space, does the equality $\text{DT}_{A+B_{<}(0,r)} = \text{DT}_A - r$ hold in A if we use the distance defined by the norm? If the statement is not always true, find some hypotheses under which it is true. What if we only know that we are in an abelian group with a translation-invariant metric??...

5.3. Let X be a metric space with metric d . Define a metric d_α in $X \times \mathbf{R}$ by

$$d_\alpha((x, s), (y, t)) = \alpha d(x, y) + |s - t|, \quad (x, s), (y, t) \in X \times \mathbf{R},$$

where α is any positive number. Prove that, for any subset A of X ,

$$\text{DT}_{\text{epi DT}_A}(x, s) = \text{DT}_A(x) - s, \quad (x, s) \in X \times \mathbf{R},$$

if $\alpha \geq 1$. Show by examples that this is not necessarily so if $0 < \alpha < 1$.

5.4. *Square pixels.* The centers of square pixels can be placed at the points \mathbf{Z}^2 with integer coordinates in \mathbf{R}^2 . Since \mathbf{Z}^2 is group, we can use infimal convolution. Actually we can work indifferently in \mathbf{R}^2 or \mathbf{Z}^2 . Define a function $F: \mathbf{R}^2 \rightarrow [0, +\infty]$ by $F(0) = 0$, $F(\pm 1, 0) = F(0, \pm 1) = a > 0$, $F(\pm 1, \pm 1) = b > 0$ (four points with value b), and $F(x) = +\infty$ at all other points. Then let $f = \lim F \square F \square \cdots \square F$ as the number of terms tends to infinity.

(a) Show that if $a = 1$, $b = \sqrt{2}$, then

$$1 \leq \frac{f(x)}{\|x\|_2} \leq \sqrt{\frac{2\sqrt{2}}{\sqrt{2}+1}} \approx 1.08239, \quad x \in \mathbf{Z}^2.$$

(b) Show that the best choice for a and b if we want to have

$$1 \leq \frac{f(x)}{\|x\|_2} \leq C, \quad x \in \mathbf{Z}^2,$$

with C as small as possible is $a = 1$, $b = \sqrt{2}$.

(c) For all $a, b > 0$ there is a constant C depending on a, b such that

$$C^{-1} \leq \frac{f(x)}{\|x\|_2} \leq C, \quad x \in \mathbf{Z}^2.$$

Take $b = a\sqrt{2}$ and find the a which renders C as small as possible. Show that the smallest possible value of C is

$$C = \Lambda(f, \|\cdot\|_2) = \sqrt[4]{\frac{4}{2+\sqrt{2}}} \approx 1.04038,$$

and that it is attained when a and $b = a\sqrt{2}$ take the values (5.18).

(d) Now vary both a and b . Show that the smallest C such that $C^{-1} \leq f(x)/\|x\|_2 \leq C$ for all $x \in \mathbf{Z}^2$ is the same as in (c), i.e., that we do not gain anything by taking b different from $b = a\sqrt{2}$.

(e) Prove that if we use $\Lambda'(f, \|\cdot\|_2)$ to measure deviation from the Euclidean distance, then the optimal values are as indicated in (5.19).

(f) Prove that if we use $\Lambda'(\|\cdot\|_2, f)$ to measure deviation from the Euclidean distance, then the optimal values are as indicated in (5.20).

5.5. *Hexagonal pixels.* The centers of hexagonal pixels form a triangular pattern: if we identify \mathbf{R}^2 with \mathbf{C} for convenience they can be placed at the points $p + q\omega$, $p, q \in \mathbf{Z}$, where $\omega = \frac{1}{2} + i\frac{\sqrt{3}}{2}$.

(a) Define a function $F: \mathbf{C} \rightarrow [0, +\infty]$ by $F(0) = 0$, $F(\pm 1) = F(\pm\omega) = F(\pm(1-\omega)) = 1$, and $F(x) = +\infty$ at all other points. Then let $f = \lim F \square F \square \cdots \square F$ as the number of terms tends to infinity. Show that

$$1 \leq \frac{f(x)}{\|x\|_2} \leq \frac{2}{\sqrt{3}} \approx 1.15470, \quad x \in \mathbf{Z} + \mathbf{Z}\omega.$$

Thus a suitable multiple f_1 of f has $\Lambda_2(f_1) \leq \sqrt{2}/\sqrt[4]{3} \approx 1.07457$.

(b) Now define $F: \mathbf{Z}^2 \rightarrow [0, +\infty]$ by $F(0) = 0$, $F(\pm 1, 0) = F(0, \pm 1) = F(1, -1) = F(-1, 1) = 1$ and $F(x) = +\infty$ at all other points. Let $f = \lim F \square F \square \cdots \square F$ as the number of terms tends to infinity. Find a Euclidean norm $\|\cdot\|$, i.e., a norm defined by an inner product, such that

$$\|x\| \leq f(x) \leq C\|x\|, \quad x \in \mathbf{Z}^2,$$

with C as small as possible.

5.6. *Triangular pixels.* The centers of triangular pixels form a hexagonal pattern. These hexagonally placed points do not form a subgroup of \mathbf{R}^2 , and therefore infimal convolution cannot be applied directly as in the case of square or hexagonal pixels. However, if we take two steps in the hexagons, we get a group of triangularly placed points, which is isometric to the group $\mathbf{Z} + \mathbf{Z}\omega$ formed by the hexagonal pixels considered in exercise 5.5. For an even number of steps, the inequalities for hexagonal pixels can be applied; for an odd number of steps, we first take an even number of steps and then one extra step. Investigate what the inequalities of exercise 5.5 yield.

5.7. Prove that in an abelian group with a translation-invariant metric which is upper regular for the triangle inequality, $B_{\leq}(c, r) = B_{\leq}(c', r')$ implies that $c = c'$.

5.8. Prove that the chamfer distance in \mathbf{Z}^2 defined by $F(\pm 1, 0) = F(0, \pm 1) = 3$, $F(\pm 1, \pm 1) = 4$ is not lower regular for the triangle inequality.

6. Skeletonizing

6.1. Definition of the skeleton

If A is any subset of a metric space X , then its interior A° is the union of all open balls contained in A . This is typically the union of a very large family of sets. We would like to describe A° as the union of a smaller family. It is obvious that if we have two balls contained in A , $B_{<}(c, r)$ and $B_{<}(c', r')$, and one is contained in the other, then we may throw away the smaller ball without changing the union. In fact, for every ball $B_{<}(c, r)$ in the union, we may throw away all balls contained in that ball without changing the union. This leads to the concept of a maximal ball. A maximal ball must be retained, but all balls contained in a maximal ball may be dispensed with.

The importance of skeletons in applications is due to the fact that they are thin in some sense but nevertheless retain important information about an object, for instance its general shape, and that, given the skeleton and the distance transform at the points in the skeleton, we can reconstruct the whole object. Typically we save memory when listing only the skeleton and the quench function.

If c is the center of a maximal open ball $B_{<}(c, r)$ contained in a set A , then necessarily $r = -DT_A(c)$. In fact, when we defined the distance transform $DT_A(c)$ at a point c , we looked at all balls with center c contained in A and we took the largest such ball. Note that then we kept the center fixed. There is a largest ball with center c , which in particular is maximal among these balls. By way of contrast, when we define the skeleton we shall vary both the center and the radius and look at all balls contained in A , regardless of their centers. We shall now give a name to the centers of maximal balls.

If A is a subset of a metric space X , then we shall define the *skeleton*¹⁶ of A , denoted by $\text{Sk}(A)$, as the set of all centers of maximal open balls contained in A . This means that $c \in \text{Sk}(A)$ if and only if there exists a number $r > 0$ such that $B_{<}(c, r) \subset A$ and such that if a ball $B_{<}(c', r')$ is contained in A and contains $B_{<}(c, r)$, then $B_{<}(c', r') = B_{<}(c, r)$.¹⁷ The skeleton may be empty: think of a set with empty interior or of a half-space in \mathbf{R}^n . A half-space contains lots of balls, but there are

¹⁶It seems that the term *skeleton* was first used in this context by Rosenfeld & Pfaltz (1966).

¹⁷I did not say that $c' = c$ and $r' = r$.

no maximal balls. So obviously we need to investigate whether there exist maximal balls—and whether there are enough of them in the formation of the interior of A . For this we shall need Zorn's Lemma.

6.2. Existence of skeletons

An ordered set X is said to be *totally ordered* if for any two elements $x, y \in X$ we have $x \leq y$ or $y \leq x$. An ordered set X is said to be *inductive* or *inductively ordered* (Bourbaki (1963:34)) if every totally ordered subset of X possesses a majorant in X .¹⁸ This means that for every $Y \subset X$ which is totally ordered, there exists an element $b \in X$ such that $y \leq b$ for all $y \in Y$. This concept is of interest because it is used as an hypothesis in Zorn's Lemma, which guarantees the existence of maximal elements.

Theorem 6.1 (Zorn's Lemma). *Every inductively ordered set possesses a maximal element.*

We shall accept Zorn's Lemma here. One can prove it using the Axiom of Choice; conversely, the latter can be proved from Zorn's Lemma. To establish the existence of a maximal element we shall have to prove that a certain order is inductive.

Theorem 6.2. *Let \mathbf{Z}^n be equipped with a metric which either is inherited from a norm on \mathbf{R}^n or a chamfer distance in the sense of Definition 5.18, and let A be a finite subset. Then the set of all open balls contained in A is inductively ordered.*

Proof. Let us consider a union $A_M = \bigcup_{(c,r) \in M} B_{<}(c, r)$ of a family of open balls, where M is a subset of $\mathbf{Z}^n \times \mathbf{R}$. Assume that the family is totally ordered, i.e., for any two pairs $(c, r), (c', r') \in M$, either $B_{<}(c, r)$ is a subset of $B_{<}(c', r')$ or conversely. The set of radii must be bounded since A is finite. So clearly A_M is finite, which implies that it is equal to one of the balls $B_{<}(c, r)$ with $(c, r) \in M$. We are done.

In \mathbf{R}^n things are less simple.

Theorem 6.3. *Let A be a set in a finite-dimensional normed vector space E . Assume that A is bounded. Then the set of open balls contained in A is inductively ordered.*

Proof. Let as before $A_M = \bigcup_{(c,r) \in M} B_{<}(c, r)$ be a union of open balls, where M is a subset of $E \times \mathbf{R}$. We assume that the family of these balls is totally ordered. Define R as the supremum of all numbers r such that there exists a $c \in E$ such that $(c, r) \in M$. Since A is bounded, this supremum is necessarily finite. For every $j = 1, 2, 3, \dots$ there exists a number $r_j > R - 1/j$ and a point c_j in A such that $(c_j, r_j) \in M$. Unless one of the radii is equal to R we may also choose the r_j so that $r_{j+1} > r_j$. Since the sequence of centers (c_j) is bounded, it has a converging subsequence; let us change notation so that (c_j) itself is converging. Let its limit be C . We claim that A_M is equal to the ball $B_{<}(C, R)$.

We shall prove first that A_M is contained in $B_{<}(C, R)$. Let $x \in A_M$; there exists $(c, r) \in M$ such that $d(x, c) < r$. Define $\varepsilon = r - d(x, c) > 0$. Then the ball $B_{<}(x, \varepsilon)$ is contained in $B_{<}(c, r)$. If $r = R$, then we must have $B_{<}(c, r) = B_{<}(C, R)$. If, on the other hand $r < R$, then we take k so large that $r_k > r$, and we must have $B_{<}(c, r) \subset B_{<}(c_k, r_k)$ since the opposite inclusion is impossible and one of them must hold by hypothesis. Thus $B_{<}(x, \varepsilon) \subset B_{<}(c_j, r_j)$ for all $j \geq k$, which implies, on letting

¹⁸The empty set is totally ordered, so an inductively ordered set must be nonempty.

j tend to infinity, that $B_{<}(x, \varepsilon) \subset B_{\leq}(C, R)$. But then x belongs also to the open ball $B_{<}(C, R)$ and we are done. (Note that this part of the proof is not valid in a general metric space.)

In the other direction, let us prove that $B_{<}(C, R)$ is contained in A_M . Take any point $x \in B_{<}(C, R)$. Then $d(x, C) < R$, and we may define $\varepsilon = R - d(x, C) > 0$. We then have $d(x, c_j) \leq d(x, C) + d(C, c_j) = R - \varepsilon + d(C, c_j)$, where the right-hand side is less than r_j for large j . Hence x belongs to $B_{<}(c_j, r_j)$, and therefore to A_M . (This part of the proof is valid in any compact metric space.)

If the norm is Euclidean, it is enough to assume that A does not contain a half-space. Also, for any given norm in \mathbf{R}^n , it is enough to assume that A does not contain a cone of a certain aperture.

Example. One might think that the result should hold in any compact metric space. However, simple examples show that this is not so. Define a compact metric space X as consisting of the segment $A = [-1, 1] \times \{0\}$ in \mathbf{R}^2 and join to it the point $(0, 1)$. The metric shall be that induced by the usual Euclidean metric in \mathbf{R}^2 . Consider now the open balls $B_{<}(c_j, r_j)$, where $c_j = (1/j, 0)$ and $r_j = 1$, $j = 1, 2, \dots$. Note that $B_{<}(c_j, r_j) =]-1 + 1/j, 1] \subset A$. The sequence of balls is increasing with j , and its union is the segment $A_M =]-1, 1] \times \{0\}$. However, this segment is not an open ball in X . What is worse: it is not contained in an open ball in A . Hence the family of open balls in A is not inductively ordered.

There is a well-known trick to get rid of the hypothesis that A must not contain any half-space. Indeed a half-space is a limiting case of a ball, and if we compactify \mathbf{R}^n by adding a point at infinity, then $\mathbf{R}^n \cup \{\infty\}$ can be regarded as the n -sphere. The balls in \mathbf{R}^n then become spherical caps on the sphere, and so do the half-spaces. But these spherical caps are actually balls for a suitable metric on the sphere. So $\mathbf{R}^n \cup \{\infty\}$ is a metric space and the open balls contained in any given subset are inductively ordered—the proof is very much like the one we have done in \mathbf{R}^n . The compactification allows us to define a generalized skeleton, which may contain ∞ , but we need of course to modify the quench function.

Corollary 6.4. *Let A be a bounded subset of a finite-dimensional normed vector space, or a bounded subset of \mathbf{Z}^n , where \mathbf{Z}^n is provided with a metric as in Theorem 6.2. The union of all open balls with center c belonging to the skeleton and radius equal to $-DT_A(c)$ is equal to the interior of A . In particular, if A has interior points, then the skeleton of A is nonempty.*

Proof. Take any point $x \in A^\circ$. The ball $B_{<}(x, \varepsilon)$ is contained in A° for some small positive ε . By Zorn's lemma (Theorem 6.1) and Theorem 6.2 or 6.3, respectively, there is a maximal ball $B_{<}(c, r)$ containing $B_{<}(x, \varepsilon)$ and contained in A . Thus $c \in \text{Sk}(A)$ and $x \in B_{<}(c, r)$, with $r = -DT_A(c)$.

In any metric space where the conclusion of Theorem 6.3 holds we have $A^\circ = \bigcup_{c \in \text{Sk}(A)} B_{<}(c, -DT_A(c))$. Here $-DT_A(c) = (DT_A(c))^-$ is the quench function evaluated at c . Knowledge of $\text{Sk}(A)$ and the restriction of DT_A to $\text{Sk}(A)$ is equivalent to knowing A° . This shows how we can reconstruct A° from $\text{Sk}(A)$ and the quench function.

6.3. Properties of skeletons

In some sense the skeleton is a thin set. For instance, it is easy to prove that the skeleton in \mathbf{R}^n has no interior points (cf. exercises 6.1 and 6.2). On the other hand, the closure of the skeleton need not be of Lebesgue measure zero (cf. exercise 6.3). These results are mentioned in Serra (1982:378) and Matheron (1988:218). It is also stated there that it is unknown whether the skeleton has Lebesgue measure zero, and whether the interior of its closure is empty.¹⁹

The skeleton of a set A in \mathbf{R}^2 need not be a closed set, even if A has a smooth boundary. To construct an example, let first $U = A + B$, the dilation of the segment $A = [(-1, 1), (1, 1)]$ by the unit disk $B = \{(x, y) \in \mathbf{R}^2; x^2 + y^2 < 1\}$. This is an open set and its skeleton is $\text{Sk}(U) = A$. Then modify U as follows. The lower boundary of U is just the segment $[(-1, 0), (1, 0)]$, so that U is locally defined by the inequality $y > 0$ near this segment. We replace it by an inequality $y > \varphi(x)$, where $\varphi(x) = \sum_1^\infty c_j \psi(2^j x - 1)$, ψ being an even, nonpositive function in C_0^∞ with $\psi(0) = -1$ and support contained in $[-\frac{1}{3}, \frac{1}{3}]$. For a suitable choice of positive constants c_j , the function φ is in C^∞ , and the skeleton of the new open set V so defined contains segments $\{2^{-j}\} \times [a, b]$ for $j = 1, 2, 3, \dots$ for a suitable choice of a and $b > a$. But no point on the limiting segment $\{0\} \times [a, b]$, belongs to it. (This example is essentially taken from Matheron (1988:219).)

The following result was proved in \mathbf{R}^n by Matheron (1988:225).

Theorem 6.5. *Let E be a normed space with metric given by the norm: $d(x, y) = \|x - y\|$. Let A be a nonempty proper subset of E , fix a point c in the interior of A , and define $h(x) = d(x, c) + \text{DT}_A(x)$, $x \in E$. Then c belongs to the skeleton of A if and only if h has a minimum only at c .*

Proof. If $B_<(c, r) \subset B_<(x, r') \subset A$, where $r = -\text{DT}_A(c) > 0$, then $r' \geq d(c, x) + r$. This implies that $h(x) = d(x, c) + \text{DT}_A(x) \leq d(x, c) - r' \leq d(c, c) - r = h(c)$. If c is the only point where h attains its minimum, we must have $x = c$ and it follows that $B_<(c, r)$ is maximal, hence that $c \in \text{Sk}(A)$.

Conversely, assume that c is in the skeleton and that x is a point where $h(x) \leq h(c)$. Then $d(c, x) - r' \leq -r$, where we define $r = -\text{DT}_A(c)$, $r' = -\text{DT}_A(x)$. This implies that $B_<(x, r')$ contains $B_<(c, r)$. Since c is in the skeleton of A , the two balls must be equal, which implies that $x = c$. Therefore the infimum of h is attained at c and only there.

Thanks to the calculus of balls developed in section 5.6 we can generalize this result to other groups. In a normed space of positive dimension, the open ball of radius $r = -\text{DT}_A(c)$ is the interior of the closed ball of the same radius and the same center. In a group where the set of distances is discrete, the open ball $B_<(c, r)$ can be described as the closed ball of radius $\rho_-(r)$. Since the conditions for working with closed balls are more easily satisfied than those for open balls, we will get a more applicable result if we replace the function $x \mapsto d(x, c) + \text{DT}_A(x)$ by $x \mapsto d(x, c) - \rho_-(-\text{DT}_A(x))$.

Theorem 6.6. *Let G be an abelian group with a translation-invariant metric d which is upper regular for the triangle inequality and such that the set of all distances is discrete. Let A be a nonempty proper subset of G , fix a point $c \in A$, and define*

¹⁹It seems that the answers to these questions are still unknown.

$h(x) = d(x, c) - \rho_-(-DT_A(x))$, $x \in G$. Then c belongs to the skeleton of A if and only if h has a minimum only at c .

Proof. If $B_{<}(c, r) \subset B_{<}(x, r') \subset A$, where $r = -DT_A(c)$, then by Theorem 5.22 $\rho_-(r') \geq d(c, x) + \rho_-(r)$. This implies that $h(x) = d(x, c) - \rho_-(-DT_A(x)) \leq d(x, c) - \rho_-(r') \leq d(c, c) - \rho_-(r) = h(c)$. If c is the only point where h attains its minimum, we must have $x = c$ and it follows that $B_{<}(c, r)$ is maximal, hence that $c \in \text{Sk}(A)$.

Conversely, assume that c is in the skeleton and that x is a point where $h(x) \leq h(c)$. Then $d(c, x) - \rho_-(r') \leq -\rho_-(r)$, where we define $r = -DT_A(c)$, $r' = -DT_A(x)$. This implies that $B_{\leq}(x, \rho_-(r'))$ contains $B_{\leq}(c, \rho_-(r))$. But $B_{\leq}(x, \rho_-(r')) = B_{<}(c, r')$ and $B_{\leq}(c, \rho_-(r)) = B_{<}(c, r)$ by Proposition 5.21. Since c is in the skeleton of A , the two balls must be equal, which implies that the difference between their centers belongs to a bounded subgroup, thus that $x = c$ (cf. exercise 5.7). Therefore the infimum of h is attained at c and only there.

Continuity properties of skeletons....

6.4. Exercises

6.1. Show that in \mathbf{R}^n with a Euclidean metric every skeleton has an empty interior.

6.2. Let G be an abelian group and P any nonempty subset of G . A set $A \subset G$ will be called P -open if it is of the form $\bigcup_{x \in M} (x + P)$ for some $M \subset G$. The set of all P -open sets form what we may call a semitopology τ_P on G : an arbitrary union of P -open sets is P -open. Assume now that P is finite, symmetric and contains the origin, let F be an even function satisfying the conditions in Corollary 5.17, and let f obtained from F as in that corollary. Let finally $\text{Sk}(A)$ be the skeleton of a set A defined by the distance $d(x, y) = f(x - y)$. Prove that interior of $\text{Sk}(A)$ defined by τ_P is empty: there is no nonempty P -open set contained in $\text{Sk}(A)$.

6.3. Construct an example of an open set in \mathbf{R}^2 such that the closure of its skeleton with respect to a Euclidean metric is of positive Lebesgue measure. *Hint:* Take U as the unit disk, and add to it denumerably many equilateral triangles with two points on the circumference and the third point with rational argument. This will yield a set whose skeleton contains a lot of rays; the limits of these rays have arguments in a set of measure $2\pi - \varepsilon$.

6.4. Show that if we use the l^∞ distance there exists a bounded open connected set in \mathbf{R}^2 such that its skeleton is not connected.

6.5. We define the r -skeleton, associated to the radius $r > 0$, as the set $\text{Sk}_r(A)$ of points in the skeleton where the quench function is equal to r . Therefore $\text{Sk}(A) = \bigcup_{r>0} \text{Sk}_r(A)$.

Assume now that we are in a normed vector space with distance $d(x, y) = \|x - y\|$.

(a) Find some conditions on a set A which guarantee that the skeleton of its dilation by a ball is the same as that of A :....

$$\text{Sk}(A + B_{<}(0, R)) = \text{Sk}(A).$$

(b) Find conditions on a set A which ensure that the skeleton of an erosion is simply obtained by deleting a certain part of the skeleton:....

$$\text{Sk}(A \ominus B_{<}(0, R)) = \bigcup_{r>R} \text{Sk}_r(A).$$

(c) Show that

$$\text{Sk}_R(A) = \bigcap_{r>0} \left[(A \ominus B_{<}(0, R)) \setminus (A \ominus B_{<}(0, R))_{B_{\leq}(0, r)} \right].$$

(d) Show that the reconstruction of the interior of A takes the form??.... Hjälper det om A inte innehåller ett halvrum??....

$$A^\circ = \bigcup_{r>0} (\text{Sk}_r(A) + B_{<}(0, r)).$$

(d) Show that we can also reconstruct the dilations and erosions by balls from the r -skeletons:??....

$$A + B_{<}(0, R) = \bigcup_{r>0} (\text{Sk}_r(A) + B_{<}(0, R + r)),$$

$$A \ominus B_{<}(0, R) = \bigcup_{r>R} (\text{Sk}_r(A) + B_{<}(0, r - R)).$$

6.6. What remains of exercise 6.5 if we only know that we are in an abelian group with a translation-invariant metric?

7. Lattices

7.1. Definition and first properties of lattices

Lattice theory was developed by Garrett Birkhoff and others in the beginning of the twentieth century (Birkhoff 1940, 1948). There is an analogy between lattice theory and the theory of vector spaces. The theory of topological vector spaces was developed to a large extent because of the theory of distributions, which in turn was motivated by applications in partial differential equations. Developments in image processing motivated a renewed interest in lattice theory, in particular in complete lattices. Lattice theory was applied to switching circuits, and it was then enough, because of general finiteness conditions, to form models using lattices, but in image processing it is more convenient to assume completeness.

While vector spaces are useful in modelling linear problems, lattices seem to be more adapted to nonlinear problems. Auditory phenomena are often additive: all the instruments of an orchestra can be heard; while with visual phenomena this is not so: one object can block another from our view. This indicates that linear models may suffice for the first kind of phenomena, while the visual ones are more in agreement with nonlinear operators like supremum and infimum.

There are also analogies between preordered sets, in particular lattices, and topological spaces as we shall see. The increasing mappings in the first case correspond to continuous mappings in the second.

Let L be an ordered set and A a subset of L . An element $b \in L$ is said to be the *infimum* of all elements $a \in A$ if b is the largest minorant of all $a \in A$. This means that $b \leq a$ for all elements $a \in A$, and that if $b' \leq a$ for all $a \in A$, then $b' \leq b$. The infimum, if it exists, is necessarily unique. The infimum over the empty set exists if and only if L possesses a largest element, and if so the infimum is this largest element. Similarly $\sup_{x \in \emptyset} x$ exists if and only if L has a smallest element.

We shall write $b = \inf_{a \in A} a = \inf(a; a \in A) = \bigwedge_{a \in A} a$; if A has only n elements we write $b = a_1 \wedge \cdots \wedge a_n$, in particular $b = a_1 \wedge a_2$ if $n = 2$. Similarly we define the supremum $c = \sup_{a \in A} a = \bigvee_{a \in A} a$ as the smallest majorant of all elements in A . If any set consisting of two elements in L has an infimum, we shall call L an *inf-semilattice*; similarly, if any two-set of L has a supremum, we shall call L a *sup-semilattice*. If L is both an inf-semilattice and a sup-semilattice we shall call L a *lattice*.

If any nonempty subset, finite or infinite, has an infimum, L will be said to be a *complete inf-semilattice*; analogously we define *complete sup-semilattice* and *complete lattice*. A complete inf-semilattice has a smallest element, which is the infimum of all elements, denoted by $\mathbf{0}$, and a complete sup-semilattice has a largest element, the supremum of all elements, denoted by $\mathbf{1}$.

A complete inf-semilattice with a largest element $\mathbf{1}$ is also a complete lattice. Indeed, the supremum of any set of elements is equal to the infimum of all majorants of the set—this set is not empty since $\mathbf{1}$ is a majorant.

In a complete lattice, the infimum of the empty set exists and is $\mathbf{1}$, and the supremum over the empty set is $\sup_{j \in \emptyset} x_j = \mathbf{0}$.

It is possible to define a lattice as a set with two binary operations \wedge and \vee satisfying certain axioms.

A *sublattice* is defined just like a subgroup with respect to the operations \wedge and \vee : that M is a sublattice of L means that for all $x, y \in M$, $x \wedge y$ and $x \vee y$, when calculated in L , are elements of M . A sublattice is therefore something more than a subset with the induced order; see the following examples.

Example. The space of real-valued continuous functions on a topological space is a lattice with the usual order: $f \leq g$ if and only if $f(x) \leq g(x)$ for all x . The space $C^1(\mathbf{R}^n)$ of continuously differentiable functions on \mathbf{R}^n is not a sublattice of $C(\mathbf{R}^n)$. It is not even a lattice on its own.

Example. The set $[-\infty, +\infty]^{\mathbf{R}^n}$ of all functions defined on \mathbf{R}^n and with values in the extended real line is a lattice under the usual order for real numbers, extended in an obvious way to the two infinities. The subset of all convex functions is ordered in the same way, and is also a lattice under this order. However, the convex functions $CVX(\mathbf{R}^n)$ do not form a sublattice of $[-\infty, +\infty]^{\mathbf{R}^n}$. The supremum of two convex functions is equal to the pointwise supremum of them:

$$f \vee g = \max(f, g),$$

but the infima are different in the two lattices: the infimum in the lattice of convex functions is

$$f \wedge_{cvx} g = \sup [h \in CVX(\mathbf{R}^n); h \leq f, g] \leq \min(f, g),$$

where the supremum is calculated in $[-\infty, +\infty]^{\mathbf{R}^n}$ and has a sense because that lattice is complete.

Example. The family $\mathcal{P}(W)$ of all subsets of a set W is a complete lattice, with $\bigwedge A_j = \bigcap A_j$ and $\bigvee A_j = \bigcup A_j$. The compact sets in \mathbf{R}^n form a sublattice $\mathcal{C}(\mathbf{R}^n)$ of $\mathcal{P}(\mathbf{R}^n)$. This lattice is a complete inf-semilattice but not a complete sup-semilattice. The family $\mathcal{K}(\mathbf{R}^n)$ of all convex compact sets is a lattice but not a sublattice of $\mathcal{C}(\mathbf{R}^n)$: the supremum of two compact sets is not the same in the two lattices.

Example. The family of all closed sets in \mathbf{R}^n , denoted by $\mathcal{F}(\mathbf{R}^n)$, is a sublattice of $\mathcal{P}(\mathbf{R}^n)$: the union and intersection of two closed sets are closed. But, although $\mathcal{F}(\mathbf{R}^n)$ is a complete lattice, it is not a sub-complete-lattice of the complete lattice $\mathcal{P}(\mathbf{R}^n)$. The union of a family of closed sets is not always closed, but there is a supremum, viz. the closure of the union. Thus, finite suprema agree with those in $\mathcal{P}(\mathbf{R}^n)$ while infinite suprema do not. This example shows that we would need a different word for *complete lattice*, to allow for a better term than “sub-complete-lattice.”

7.2. Morphology on lattices

We already defined the epigraph of a mapping $X \rightarrow [-\infty, +\infty]$; see (2.8). The definition makes of course sense for any mapping $X \rightarrow Y$, where X is an arbitrary set and Y a preordered set, thus

$$(7.1) \quad \text{epi } f = \{(x, y) \in X \times Y; f(x) \leq y\}.$$

It will be convenient to define also the *hypograph* of a mapping $f: X \rightarrow Y$; it is

$$(7.2) \quad \text{hypo } f = \{(x, y) \in X \times Y; y \leq f(x)\}.$$

The intersection of the two, $\text{epi } f \cap \text{hypo } f$, contains the *graph* of f , which is the set of all pairs $(x, f(x))$:

$$(7.3) \quad \text{graph } f = \{(x, y) \in X \times Y; y = f(x)\}.$$

If Y is not only preordered but ordered, the intersection is equal to the graph. If Y has the discrete order, then $\text{epi } f = \text{hypo } f = \text{graph } f$; if Y has the chaotic preorder, then $\text{epi } f = \text{hypo } f = X \times Y$.

In $L' = \text{Incr}(L, L)$, the set of all increasing mappings of an ordered set into itself, the idempotent mappings are of particular interest. We shall call an increasing and idempotent mapping from L into L a *morphological filter*, and denote the set of all such mappings by $\text{Filt}(L)$. Openings and closings are examples: an opening is an antiextensive mapping in $\text{Filt}(L)$, and a closing is an extensive mapping in $\text{Filt}(L)$.

We have studied dilations in an abelian group, i.e., mappings of the form $A \mapsto A+B$. We also remarked (cf. (2.11)) that a mapping which commutes with translations and the formation of infinite unions is necessarily of this form. In lattice theory it is therefore natural to say that a mapping $\delta: L \rightarrow M$, where L and M are complete lattices, is a *dilation* if it commutes with the formation of suprema, i.e., $\delta(\bigvee_j x_j) = \bigvee_j \delta(x_j)$ for all indexed families $(x_j)_{j \in J}$. In particular we get $\delta(\mathbf{0}_L) = \mathbf{0}_M$ (take J empty), while $\delta(\mathbf{1}_L) = \bigvee_{x \in L} \delta(x) \leq \mathbf{1}_M$ (give an example where $\delta(\mathbf{1}_L) \neq \mathbf{1}_M$). Similarly we shall say that ε is an *erosion* if it commutes with the formation of infinite intersections, $\varepsilon(\bigwedge_j x_j) = \bigwedge_j \varepsilon(x_j)$. We note that $\varepsilon(\mathbf{1}_L) = \mathbf{1}_M$ (take J empty), while $\varepsilon(\mathbf{0}_L) = \bigwedge_{x \in L} \varepsilon(x) \geq \mathbf{0}_M$.

Dilations and erosions are always increasing. Indeed, we have $\delta(x \vee y) = \delta(x) \vee \delta(y)$. If $x \leq_L y$, this equation simplifies to $\delta(y) = \delta(x) \vee \delta(y) \geq_M \delta(x)$, which shows that δ is increasing, $\delta \in \text{Incr}(L, M)$. A similar argument shows that erosions are increasing.

We now recall the exercise 2.2(a) where it was stated that $A + B \subset C$ in an abelian group if and only if $A \subset C \ominus \check{B}$; in other words that $\delta_B(A) \subset C$ if and only if $A \subset \varepsilon_{\check{B}}(C)$. In a lattice this may be written as $\delta(x) \leq y$ iff $x \leq \varepsilon(y)$, equivalently as

$\text{epi } \delta = (\text{hypo } \varepsilon)^{-1}$, where the exponent -1 means that we swap the components: for a subset A of a Cartesian product $X \times Y$ we define $A^{-1} = \{(y, x); (x, y) \in A\} \subset Y \times X$. May we use this as a model to define erosions from dilations and conversely in the more general lattice situation? Indeed this is the case, and we shall now show this.

7.3. Inverses of mappings between lattices

Let L be a complete lattice, M a preordered set, and $f: L \rightarrow M$ any mapping. We then define the *upper inverse* $f^{[-1]}: M \rightarrow L$ and the *lower inverse* $f_{[-1]}: M \rightarrow L$ as the mappings

$$(7.4) \quad f^{[-1]}(y) = \bigwedge_{x \in L} (x; f(x) \geq_M y) = \bigwedge_{x \in L} (x; (x, y) \in \text{hypo } f), \quad y \in M;$$

$$(7.5) \quad f_{[-1]}(y) = \bigvee_{x \in L} (x; f(x) \leq_M y) = \bigvee_{x \in L} (x; (x, y) \in \text{epi } f), \quad y \in M.$$

As a first observation, let us note that these inverses are always increasing. If there exists a smallest element $\mathbf{0}_M$, then $f^{[-1]}(\mathbf{0}_M) = \mathbf{0}_L$. Similarly, if M possesses a largest element $\mathbf{1}_M$, then $f_{[-1]}(\mathbf{1}_M) = \mathbf{1}_L$. If M has the chaotic preorder, then both inverses are constant, $f^{[-1]} = \mathbf{0}_L$ and $f_{[-1]} = \mathbf{1}_L$ identically.

If, given a mapping $f: L \rightarrow M$, we can find a mapping $g: M \rightarrow L$ such that $\text{epi } g = (\text{hypo } f)^{-1}$ we would be content to have a kind of inverse to f . However, usually the best we can do is to study mappings with $\text{epi } g \supset (\text{hypo } f)^{-1}$ or $\text{epi } h \subset (\text{hypo } f)^{-1}$. This we shall do in the following proposition, which shows that the upper and lower inverses are solutions to certain extremal problems.

Proposition 7.1. *Let L be a complete lattice, M a preordered set, and let $f: L \rightarrow M$, $g, h: M \rightarrow L$ be mappings. If $\text{epi } g \supset (\text{hypo } f)^{-1} \supset \text{epi } h$, then $g \leq f^{[-1]} \leq h$ and*

$$\text{epi } g \supset \text{epi } f^{[-1]} \supset (\text{hypo } f)^{-1} \supset \text{epi } h.$$

Hence $f^{[-1]}$ is the largest mapping g such that $\text{epi } g \supset (\text{hypo } f)^{-1}$. If on the other hand $\text{hypo } g \subset (\text{epi } f)^{-1} \subset \text{hypo } h$, then $g \leq f_{[-1]} \leq h$ and

$$\text{hypo } g \subset (\text{epi } f)^{-1} \subset \text{hypo } f_{[-1]} \subset \text{hypo } h.$$

Hence $f_{[-1]}$ is the smallest mapping h which satisfies $\text{hypo } h \supset (\text{epi } f)^{-1}$.

The proof is straightforward.

Corollary 7.2. *With f , g and h given as in the proposition, assume that $\text{epi } g = (\text{hypo } f)^{-1}$. Then $g = f^{[-1]}$. Similarly, if $\text{hypo } h = (\text{epi } f)^{-1}$, then $h = f_{[-1]}$. If also M is a complete lattice, then $\text{epi } g = (\text{hypo } f)^{-1}$ implies that $g_{[-1]} = f$. Similarly, $\text{hypo } h = (\text{epi } f)^{-1}$ implies $h^{[-1]} = f$.*

An ideal inverse g would satisfy $g \circ f = \text{Id}_L$, $f \circ g = \text{Id}_M$, and the inverse of g would be f . It is therefore natural to compare $f^{[-1]} \circ f$ and $f_{[-1]} \circ f$ with Id_L ; $f \circ f^{[-1]}$ and $f \circ f_{[-1]}$ with Id_M ; and $(f_{[-1]})^{[-1]}$ and $(f^{[-1]})_{[-1]}$ with f . This is what we shall do now.

Proposition 7.3. *If L is a complete lattice and M a preordered set, then for all mappings $f: L \rightarrow M$ one has $f^{[-1]} \circ f \leq \text{Id}_L \leq f_{[-1]} \circ f$ with equality if f is coincreasing*

as defined in (3.5). Conversely, if $f^{[-1]} \circ f = \text{Id}_L$ or $f_{[-1]} \circ f = \text{Id}_L$, then f is coincreasing.

Proof. It is clear that $f^{[-1]}(f(a)) = \bigwedge(x; f(x) \geq f(a)) \leq a$. If f is coincreasing, then $\{x; f(x) \geq f(a)\}$ is contained in $\{x; x \geq a\}$, so that $f^{[-1]}(f(a)) \geq \bigwedge(x; x \geq a) = a$.

Conversely, if $f^{[-1]}(f(a)) \geq a$, then for all x , $f(x) \geq f(a)$ implies $x \geq a$. If this is true for all a , then f is coincreasing.

Proposition 7.4. *If g, h are increasing mappings such that $g \circ f \leq \text{Id}_L \leq h \circ f$, then $g \leq f^{[-1]}$ and $h \geq f_{[-1]}$. Hence, in view of Proposition 7.3, $f^{[-1]}$ is the largest increasing mapping g such that $g \circ f \leq \text{Id}_L$, and $f_{[-1]}$ is the smallest increasing mapping h such that $h \circ f \geq \text{Id}_L$.*

Proof. If g and h are increasing and $g \circ f \leq \text{Id}_L \leq h \circ f$, then $\text{epi } g \supset (\text{hypo } f)^{-1}$ and $\text{hypo } h \supset (\text{epi } f)^{-1}$. We can now apply Proposition 7.1.

Proposition 7.5. *Let L be a complete lattice and M a preordered set. Then $f^{[-1]}(y) \leq f_{[-1]}(y)$ for all $y \in \text{im } f$, and also for all y majorizing or minorizing $\text{im } f$. In particular, $f^{[-1]} \leq f_{[-1]}$ if f is surjective.*

Proof. If $y = f(a) \in \text{im } f$, then $f_{[-1]}(y) = \bigvee(x; f(x) \leq y) \geq a$ and $f^{[-1]}(y) = \bigwedge(x'; f(x') \geq y) \leq a$. If y majorizes all elements in $\text{im } f$, then $f_{[-1]}(y) = \mathbf{1}$, and if y minorizes all elements in $\text{im } f$, then $f^{[-1]}(y) = \mathbf{0}$.

Proposition 7.6. *Let L be a complete lattice and M a preordered set. Then the following four conditions are equivalent.*

- (a) f is coincreasing;
- (b) $f^{[-1]} \circ f \geq \text{Id}_L$;
- (c) $f_{[-1]} \circ f \leq \text{Id}_L$;
- (d) $f_{[-1]} \leq f^{[-1]}$.

Proof. (a) implies (b) and (c). If f is coincreasing we know already from Proposition 7.3 that (b) and (c) hold with equality.

(b) or (c) implies (a). If (b) or (c) holds, then, in view of Proposition 7.3, they hold with equality, and f is coincreasing.

(a) implies (d). Assume that f is coincreasing and fix an element $y \in M$. Let $x, x' \in L$ be such that $f(x) \leq y \leq f(x')$. Then $x \leq x'$. Letting x vary we see that $f_{[-1]}(y) \leq x'$. Letting x' vary, we see that $f_{[-1]}(y) \leq f^{[-1]}(y)$. Thus (d) holds.

(d) implies (a). If x and x' are given with $f(x) \leq f(x')$ we define $y = f(x)$. Then $x \leq f_{[-1]}(y)$ and $f^{[-1]}(y) \leq x'$. If we know that $f_{[-1]}(y) \leq f^{[-1]}(y)$ it follows that $x \leq x'$, proving that f is coincreasing.

Next we compose $f_{[-1]}$ with f in the other order. This will lead to a characterization of dilations—and, by duality, of erosions.

Theorem 7.7. *If L and M are complete lattices and $f: L \rightarrow M$ any mapping, then the following five properties are equivalent.*

- (A) f is a dilation;
- (B) $\text{epi } f \supset (\text{hypo } f_{[-1]})^{-1}$;

- (C) $\text{epi } f = (\text{hypo } f_{[-1]})^{-1}$;
 (D) f is increasing and $\text{epi } f \supset (\text{graph } f_{[-1]})^{-1}$;
 (E) f is increasing and $f \circ f_{[-1]} \leq \text{Id}_M$.

Proof. (A) implies (B). Suppose that (A) holds. Then if $(y, x) \in \text{hypo } f_{[-1]}$, in other words if $x \leq f_{[-1]}(y)$, we obtain, since f is increasing by hypothesis,

$$f(x) \leq f(f_{[-1]}(y)) = f(\bigvee(x; f(x) \leq y)) = \bigvee(f(x); f(x) \leq y) \leq y,$$

which means that $(x, y) \in \text{epi } f$. Thus (B) holds.

(B) implies (A). We note first that f is increasing if (B) holds. Indeed, if $x \leq x'$ and we define $y = f(x')$, then $f_{[-1]}(y) = f_{[-1]}(f(x')) \geq x' \geq x$ (see Proposition 7.3), which by (B) implies that $f(x) \leq y = f(x')$.

Let now any family $(x_j)_j$ of elements of L be given and define $z = \bigvee f(x_j)$, $w = f(\bigvee x_j)$. Since f is increasing we always have $z \leq w$. Is it true that $w \leq z$? We note that, by Proposition 7.3, $x_j \leq f_{[-1]}(f(x_j)) \leq f_{[-1]}(z)$. Taking the supremum over all j we obtain $\bigvee x_j \leq f_{[-1]}(z)$, which by (B) implies that $w = f(\bigvee x_j) \leq z$. We have proved (A).

(B) is equivalent to (C). This is clear since we always have $\text{epi } f \subset (\text{hypo } f_{[-1]})^{-1}$.

(B) implies (D). We have just seen that (B) implies that f is increasing. That (B) implies $(\text{graph } f_{[-1]})^{-1} \subset \text{epi } f$ follows from the inclusion $\text{graph } f_{[-1]} \subset \text{hypo } f_{[-1]}$.

(D) implies (B). If $x \leq f_{[-1]}(y)$ we define $x' = f_{[-1]}(y)$ and note that $x \leq x'$ and that $(y, x') \in \text{graph } f_{[-1]}$. If (D) holds, we conclude that $f(x') \leq y$. Hence, if f is increasing, $f(x) \leq f(x') \leq y$, proving (B).

(D) and (E) are equivalent: (E) is just a rewriting of (D).

Remark. We may use (B) to define dilations $L \rightarrow M$ when M is only a preordered set.

Corollary 7.8. *If L and M are complete lattices and $f: L \rightarrow M$ and $g: L \rightarrow M$ are two mappings such that $\text{epi } f = (\text{hypo } g)^{-1}$, then f is a dilation and g is an erosion, and $f_{[-1]} = g$, $g^{[-1]} = f$.*

Proof. It follows from $\text{epi } f = (\text{hypo } g)^{-1}$ that f is increasing and that $f_{[-1]} = g$, hence that (D) in the theorem holds. Since (D) is equivalent to (A), we see that f is a dilation. The rest follows by duality.

Proposition 7.9. *If L and M are complete lattices and $f: L \rightarrow M$ any mapping, then quite generally $(f_{[-1]})^{[-1]} \leq f \leq (f^{[-1]})_{[-1]}$. Equality holds at the first place if and only if f is a dilation; at the second place if and only if f is an erosion.*

Proof. We always have $\text{epi } f \subset (\text{hypo } f_{[-1]})^{-1}$, i.e., $y \geq f(a)$ implies $f_{[-1]}(y) \geq a$. This yields

$$(f_{[-1]})^{[-1]}(a) = \bigwedge(y; f_{[-1]}(y) \geq a) \leq \bigwedge(y; y \geq f(a)) = f(a).$$

If f is a dilation, then, as the last theorem shows, $\text{epi } f = (\text{hypo } f_{[-1]})^{-1}$ and equality follows.

Conversely, let us note that we always have

$$\text{epi } f \subset (\text{hypo } f_{[-1]})^{-1} \subset \text{epi } (f_{[-1]})^{[-1]}.$$

Now if $(f_{[-1]})^{[-1]} = f$, then these inclusions are equalities, and we conclude that $\text{epi } f = (\text{hypo } f_{[-1]})^{-1}$, which according to Theorem 7.7 means that f is a dilation.

Theorem 7.10. *If L and M are complete lattices and $\delta: L \rightarrow M$ a dilation, then $\delta_{[-1]}: M \rightarrow L$ is an erosion. Similarly, if $\varepsilon: L \rightarrow M$ is an erosion, then $\varepsilon^{[-1]}$ is a dilation.*

Proof. We know that $g = \delta_{[-1]}$ is increasing, so we have $g(\bigwedge y_j) \leq g(y_k)$ for all k ; hence $g(\bigwedge y_j) \leq \bigwedge g(y_k)$. We need to prove the opposite inequality, $\bigwedge g(y_k) \leq g(y)$, where $y = \bigwedge y_j$. This follows from $\delta(\bigwedge g(y_k)) \leq y$. From (E) in Theorem 7.7 we learn that $\delta(g(y_k)) \leq y_k$ for all k , which implies that $\delta(\bigwedge g(y_j)) \leq y_k$, hence that $\delta(\bigwedge g(y_j)) \leq y$ as we wanted. This completes the proof.

Corollary 7.11. *For any dilation $\delta: L \rightarrow M$ we have $\delta \circ \delta_{[-1]} \circ \delta = \delta$ and $\delta_{[-1]} \circ \delta \circ \delta_{[-1]} = \delta_{[-1]}$. Dually $\varepsilon \circ \varepsilon^{[-1]} \circ \varepsilon = \varepsilon$ and $\varepsilon^{[-1]} \circ \varepsilon \circ \varepsilon^{[-1]} = \varepsilon^{[-1]}$ for any erosion $\varepsilon: L \rightarrow M$. In particular, $\delta_{[-1]} \circ \delta$ and $\delta \circ \delta_{[-1]}$ are idempotent: $\delta_{[-1]} \circ \delta \in \text{Filt}(L)$ and $\delta \circ \delta_{[-1]} \in \text{Filt}(M)$. The first is a closing in L , the second an opening in M . Also $\varepsilon^{[-1]} \circ \varepsilon$ and $\varepsilon \circ \varepsilon^{[-1]}$ are idempotent; the first an opening, the second a closing.*

Proof. We always have $f_{[-1]} \circ f \geq \text{Id}_L$ (Proposition 7.3); it follows that $f_{[-1]} \circ f \circ f_{[-1]} \geq f_{[-1]}$. If f is increasing, we also get $f \circ f_{[-1]} \circ f \geq f$.

For dilations we have $\delta \circ \delta_{[-1]} \leq \text{Id}_M$ (Theorem 7.7), from which we deduce that $\delta \circ \delta_{[-1]} \circ \delta \leq \delta$ and $\delta_{[-1]} \circ \delta \circ \delta_{[-1]} \leq \delta_{[-1]}$. This shows what we want for dilations; the rest follows by duality.

7.4. Division of mappings between lattices

We shall now generalize the definitions of upper and lower inverses. Let three sets L, M, N be given, and let us assume that M is a complete lattice and N a preordered set. (We shall often assume that all three are complete lattices, but this is not necessary for the definitions to make sense.) Let also two mappings $f: L \rightarrow M$ and $g: L \rightarrow N$ be given. Then we may define two mappings $f/\star g, f/\star g: N \rightarrow M$ by

$$(7.6) \quad (f/\star g)(y) = \bigwedge_{x \in L} (f(x); g(x) \geq_N y), \quad y \in N,$$

$$(7.7) \quad (f/\star g)(y) = \bigvee_{x \in L} (f(x); g(x) \leq_N y), \quad y \in N.$$

We shall call them the *upper* and *lower quotient* of f and g . The quotients $f/\star g$ and $f/\star g$ increase when f increases and they decrease when g increases. Just as with division of positive numbers,

$$(7.8) \quad f_1 \leq_M f_2 \text{ and } g_1 \geq_N g_2 \text{ implies that } f_1/\star g_1 \leq_M f_2/\star g_2 \text{ and } f_1/\star g_1 \leq_M f_2/\star g_2.$$

The mappings $f/\star g$ and $f/\star g$ are always increasing. If $g(x) \geq_N y$, then $f(x) \geq_M (f/\star g)(y)$; if $g(x) \leq_N y$, then $f(x) \leq_M (f/\star g)(y)$. In particular, if $g(x) = y$, then $(f/\star g)(y) \leq_M f(x) \leq_M (f/\star g)(y)$.

If we specialize the definitions to the situation when $L = M$ and $f = \text{Id}_L$, then $f/\star g = \text{Id}_L/\star g = g^{[-1]}$ and $f/\star g = \text{Id}_L/\star g = g_{[-1]}$; see (7.4) and (7.5).

On the other hand, if we specialize to the situation where L is an arbitrary subset of a complete lattice M , $N = M$, and g is the inclusion mapping $L \rightarrow M$, then $f/\star g = f^\circ$ and $f/\star g = f_\circ$; cf. (4.6). If we specialize further, letting also f be the inclusion mapping, we obtain

$$(f/\star g)(y) = (f/\star f)(y) = f_\circ(y) = \bigvee_{x \in L} (x; x \leq y) = y^\circ \in M.$$

It is easy to verify that $y \mapsto y^\circ$ is an opening. The elements y such that $y^\circ = y$ are called L -convex. The reason should be clear from the following example.

Example. Let M be the complete lattice $[-\infty, +\infty]^E$ of functions on a vector space E , let F be a vector subspace of its dual E^\star , and let L be the set of affine functions with linear part in F , i.e., functions of the form $a(x) = \xi(x) + c$ for some linear form $\xi \in F$ and some real constant c . Then a function on E is L -convex if and only if it is equal to the supremum of all its affine minorants belonging to L . By Fenchel's theorem, this happens if and only if it is convex in the usual sense, lower semicontinuous for the topology $\sigma(E, F)$ on E generated by the linear forms in F , and not taking the value $-\infty$ except when it is equal to the constant $-\infty$.

Proposition 7.12. *If $f: L \rightarrow M$ is increasing and $g: L \rightarrow N$ is coincreasing, then $f/\star g \leq f/\star g$. If f is arbitrary and $y \in \text{img}$, then $(f/\star g)(y) \leq (f/\star g)(y)$. In particular, if g is surjective and f is arbitrary, then $f/\star g \leq f/\star g$.*

The upper quotient $f/\star g$ is the optimal solution to an inequality:

Proposition 7.13. *For all mappings $f: L \rightarrow M$ and $g: L \rightarrow N$ we have*

$$(7.9) \quad (f/\star g) \circ g \leq f \leq (f/\star g) \circ g,$$

with equality if f is increasing and g is coincreasing. Conversely, if $h, k \in \text{Incr}(N, M)$ are two increasing functions such that $h \circ g \leq f \leq k \circ g$, then $h \leq f/\star g$ and $k \geq f/\star g$. Thus $f/\star g$ is the largest increasing function h such that $h \circ g \leq f$, and $f/\star g$ is the smallest increasing function k such that $f \leq k \circ g$. In the special case $L = N$ and $g = \text{Id}_L$ we obtain

$$f/\star \text{Id}_L \leq f \leq f/\star \text{Id}_L,$$

where $f/\star \text{Id}_L$ is the largest increasing minorant of f and $f/\star \text{Id}_L$ is the smallest increasing majorant of f ; when f itself is increasing we therefore get equality.

We next compare the quotient $f/\star g$ and the composition $f \circ g^{[-1]}$:

Proposition 7.14. *For all increasing mappings $f: L \rightarrow M$ and all mappings $g: L \rightarrow N$ we have $f/\star g \geq f \circ g^{[-1]}$ with equality if f is an erosion, and $f/\star g \leq f \circ g_{[-1]}$ with equality if f is a dilation. If g is coincreasing, then $f/\star g \leq f \circ g_{[-1]} \leq f \circ g^{[-1]} \leq f/\star g$.*

Proposition 7.15. *If $h \in \text{Incr}(M, P)$, where P is an ordered set, we have $h \circ (f/\star g) \leq (h \circ f)/\star g$ with equality if h is an erosion. Similarly $h \circ (f/\star g) \geq (h \circ f)/\star g$ with equality if h is a dilation. A special case is $h \circ (f/\star \text{Id}_L) \leq (h \circ f)/\star \text{Id}_L$ (take $L = N$ and $g = \text{Id}_L$). Another special case is Proposition 7.13 (take $L = M$ and $f = \text{Id}_L$).*

Proposition 7.16. *For all mappings $f: L \rightarrow M$ we have*

$$(7.10) \quad f/\star f \leq \text{Id}_M \leq f/\star f.$$

Corollary 7.17. *For all mappings $f: L \rightarrow M$ we have*

$$(7.11) \quad (f/\star f) \circ f = f = (f/\star f) \circ f.$$

Proof. The result follows on combining (7.9), taking $g = f$, and (7.10), multiplied from the right by f —or directly from the definitions.

Theorem 7.18. *Let $f: L \rightarrow M$ be any mapping from a set L into a complete lattice M . Then $\gamma = f/\star f: M \rightarrow M$ is an opening. Conversely, any opening in M is of this form for some mapping $f: L \rightarrow M$ with $L = M$.*

Proof. It is clear that $\gamma(y) = \bigvee(f(x); f(x) \leq y)$ defines a mapping $M \rightarrow M$ which is increasing and antiextensive. The latter property implies that $\gamma \circ \gamma \leq \gamma$ (γ is an underfilter). It remains to be proved that $\gamma \circ \gamma \geq \gamma$. To do so we note that, by definition, $f(x) \leq y$ implies $f(x) \leq \gamma(y)$. Therefore

$$\gamma(y) = \bigvee(f(x); f(x) \leq y) \leq \bigvee(f(x); f(x) \leq \gamma(y)) = \gamma(\gamma(y)),$$

proving that $\gamma \leq \gamma \circ \gamma$. Note that, by (7.11), $\gamma \circ f = f$, proving that the image of f is a subset of Inv_γ ; in other words, all elements $f(x)$ are γ -open.

The last statement follows from the formula $\gamma/\star \gamma = \gamma$, which holds for any opening $\gamma: M \rightarrow M$ and is straightforward to prove. Writing out the formula in full if $\gamma = f/\star f$, we obtain $(f/\star f)/\star (f/\star f) = f/\star f$.

7.5. Exercises

7.1. Consider the following four properties of a mapping $f: X \rightarrow Y$ between two preordered sets.

- (A) f is increasing: $x \leq y \Rightarrow f(x) \leq f(y)$;
- (B) f is strictly increasing: $x < y \Rightarrow f(x) < f(y)$, where $x < y$ means that $x \leq y$ and $x \neq y$;
- (C) f is injective: $f(x) = f(y) \Rightarrow x = y$;
- (D) f is coincreasing: $f(x) \leq f(y) \Rightarrow x \leq y$.

Prove that (B) implies (A); that (A)&(C) implies (B); and that, if X is ordered, (D) implies (C).

Prove by examples that the implications (B) \Rightarrow (C); (B)&(C) \Rightarrow (D); (D) \Rightarrow (A) do not hold, not even in lattices and mappings $f: X \rightarrow X$.

If we assume (A), we thus have (D) \Rightarrow (C) and (C) \Rightarrow (B), but the converse implications do not hold.

7.2. We know that $f_{[-1]} \leq f^{[-1]}$ if f is coincreasing, i.e., satisfies (D) in the previous exercise. On the other hand $f^{[-1]} \leq f_{[-1]}$ when f is surjective. Find a lattice X and a function $f: X \rightarrow X$ such that $f^{[-1]}$ and $f_{[-1]}$ are not comparable (even with f strictly increasing and injective, i.e., satisfying (B) and (C) in the previous exercise).

7.3. (a) Formulate and prove the statement dual to Theorem 7.7.

(b) Formulate and prove the statement dual to Theorem 7.17.

(c) Formulate and prove (at least in some special cases) the statement dual to *To be, or not to be, that is the question.*

8. Notions of topology

8.1. Mappings

Let X and Y be two sets and $f: X \rightarrow Y$ a mapping from X into Y . We denote as before by $\mathcal{P}(X)$ the *power set* of X , i.e., the set of all subsets of X . We associate with f a mapping $f^*: \mathcal{P}(Y) \rightarrow \mathcal{P}(X)$ and a mapping $f_*: \mathcal{P}(X) \rightarrow \mathcal{P}(Y)$ defined by the formulas

$$(8.1) \quad f^*(B) = \{x \in X; f(x) \in B\} \text{ and } f_*(A) = \{f(x) \in Y; x \in A\}.$$

Thus $f^*(B)$, often denoted by $f^{-1}(B)$, is the *preimage* (*inverse image*) of $B \subset Y$ and $f_*(A)$, often denoted by $f(A)$, is the *direct image*, or just *image*, of $A \subset X$. It is however convenient to have a special notation for $f_*: \mathcal{P}(X) \rightarrow \mathcal{P}(Y)$, so that it is not confused with $f: X \rightarrow Y$; similarly f^* is not the pointwise inverse of f .

Given three sets X , Y and Z and a mapping $f: X \rightarrow Y$, it is customary to define a mapping $f^*: Z^Y \rightarrow Z^X$ by the formula $f^*(g) = g \circ f$, $g \in Z^Y$. The mapping $f^*(g): X \rightarrow Z$ is called the *pull-back* of g . We thus have two mappings denoted by f^* , one mapping Z^Y into Z^X as just defined, and one mapping $\mathcal{P}(Y)$ into $\mathcal{P}(X)$ defined by (8.1). Hopefully this will not cause confusion, since the two definitions are compatible in a natural way if $Z = \mathbf{R}$: if $g \in \mathbf{R}^Y$ is the characteristic function of a set $B \subset Y$, $g = \chi_B$, then $f^*(\chi_B) = \chi_{f^*(B)}$, where we use the notation from (8.1) in the right-hand side.

We note that f^* is a homomorphism of Boolean algebras: it satisfies

$$(8.2) \quad f^*(B_1 \cup B_2) = f^*(B_1) \cup f^*(B_2), \quad f^*(B_1 \cap B_2) = f^*(B_1) \cap f^*(B_2),$$

$$(8.3) \quad f^*(B_1 \setminus B_2) = f^*(B_1) \setminus f^*(B_2), \text{ in particular } f^*(\mathbb{C}B) = \mathbb{C}f^*(B).$$

The formulas (8.2) can be generalized to infinite unions and intersections, implying that f^* is both a dilation and an erosion. The homomorphism f^* is an endomorphism if and only if f is surjective, and an epimorphism if and only if f is injective.

The mapping f_* is not so well-behaved: it always satisfies

$$(8.4) \quad f_*(A_1 \cup A_2) = f_*(A_1) \cup f_*(A_2) \text{ and } f_*(\bigcup A_j) = \bigcup f_*(A_j);$$

it is thus a dilation. But for intersections we have only $f_*(A_1 \cap A_2) \subset f_*(A_1) \cap f_*(A_2)$; f_* is in general not an erosion. Also there is in general no inclusion relation between $f_*(\mathbb{C}A)$ and $\mathbb{C}f_*(A)$.

We note that $A \subset f^*(B)$ if and only if $f_*(A) \subset B$, i.e., that $(\text{hypo } f^*)^{-1} = \text{epi } f_*$. By Corollary 7.2 this implies that $(f^*)^{[-1]} = f_*$ and that $(f_*)_{[-1]} = f^*$; f_* is the dilation corresponding to the erosion f^* ; see Corollary 7.8. Moreover

$$f^* \circ f_* \circ f^* = f^* \text{ and } f_* \circ f^* \circ f_* = f_*.$$

It follows that $((f^*)^{[-1]})_{[-1]} = f^* = ((f^*)_{[-1]})^{[-1]}$ and $((f_*)^{[-1]})_{[-1]} = f_*$; cf. Proposition 7.9.

Proposition 7.3 shows that

$$(8.5) \quad f^* \circ f_* \geq \text{Id}_{\mathcal{P}(X)} \text{ and } f_* \circ f^* \leq \text{Id}_{\mathcal{P}(Y)}.$$

(These formulas can of course also be proved directly.) The first formula means that $f^*(f_*(A)) \supset A$; equality holds for all A if and only if f is injective. The second means that $f_*(f^*(B)) = B \cap \text{im} f \subset B$; equality holds for all B if and only if f is surjective. Since f^* is also a dilation, Theorem 7.7 implies that $f^* \circ (f^*)_{[-1]} \leq \text{Id}_{\mathcal{P}(X)}$. Proposition 7.3 implies that

$$(f^*)_{[-1]} \circ f^* \geq \text{Id}_{\mathcal{P}(Y)} \text{ and } (f_*)^{[-1]} \circ f_* \leq \text{Id}_{\mathcal{P}(X)}.$$

Finally we note that $(f^*)_{[-1]} = \mathbf{C} \circ f_* \circ \mathbf{C} = (f_*)^{\mathbf{d}} = ((f^*)^{[-1]})^{\mathbf{d}}$. (There seems to be no analogous formula for $(f_*)^{[-1]}$.)

We can now pass to mappings on a higher level by taking the upper and lower star again. There are four mappings on the next level:

$$(8.6) \quad (f^*)^*, (f_*)_*: \mathcal{P}(\mathcal{P}(X)) \rightarrow \mathcal{P}(\mathcal{P}(Y)) \text{ and } (f_*)^*, (f^*)_*: \mathcal{P}(\mathcal{P}(Y)) \rightarrow \mathcal{P}(\mathcal{P}(X));$$

two go in the same direction as f and two in the opposite direction. Of these four, two will be used to transport topologies, viz. $(f^*)^*$ and $(f^*)_*$.

8.2. Definition of topologies

A *topology* on a set X is a collection $\tau = \mathcal{U}(X)$ of subsets of X —thus an element of $\mathcal{P}(\mathcal{P}(X))$ —which is stable under arbitrary unions and finite intersections. The elements of $\mathcal{U}(X)$ are called *open sets*; thus any union of open sets is open and any finite intersection of open sets is open. In particular, the union and the intersection of the empty family is open, so \emptyset and X are always open subsets of X .

In a metric space there is a topology defined by the metric: it consists of all unions of balls $B_{<}(c, r)$.

A topology can be given in several different ways. We define a set as *closed* if its complement is open. Then the family $\mathcal{F}(X)$ of all closed sets is stable under arbitrary intersections and finite unions. We may also impose these conditions as axioms, and then define a set to be open if its complement is closed. A topology can be equivalently defined using open or closed sets.

Another notion is that of neighborhood. If a topology $\mathcal{U}(X)$ is given, we say that a set V is a *neighborhood* of a point x if there exists an open set U such that $x \in U \subset V$. The families $\mathcal{V}(x)$, $x \in X$, of all neighborhoods of points in X satisfy the following conditions:

$$(8.7) \quad \text{If } V \in \mathcal{V}(x), \text{ then } x \in V;$$

$$(8.8) \quad \text{If } V \in \mathcal{V}(x) \text{ and } W \supset V, \text{ then } W \in \mathcal{V}(x);$$

$$(8.9) \quad \text{If } V, W \in \mathcal{V}(x), \text{ then } V \cap W \in \mathcal{V}(x);$$

$$(8.10) \quad \text{If } V \in \mathcal{V}(x), \text{ then there exists } W \in \mathcal{V}(x) \text{ such that } V \in \mathcal{V}(y) \text{ for all } y \in W.$$

These properties are easy to verify if the topology is given and the neighborhoods are defined as above. On the other hand, if we have a collection $\mathcal{V}(x)$ for every $x \in X$ satisfying the axioms (8.7)–(8.10) and define a set U to be open if it belongs to $\mathcal{V}(x)$ for every $x \in U$, then we get a topology for which the neighborhoods are the given ones.

We can also define a topology using closure operators. If a topology is given, then we can define a closure operator in $\mathcal{P}(X)$ by taking \bar{A} as the intersection of all closed sets containing A . Then this closure operator satisfies $\overline{\bar{O}} = \bar{O}$ and $\overline{A \cup B} = \bar{A} \cup \bar{B}$. Conversely, if a closure operator is defined satisfying these two conditions we can define a set to be closed if $\bar{A} = A$; we then get a topology, a topology for which the topological closure operator is the given one.

Finally, the interior A° of a set A is the largest open set contained in the set. It is related to the closure by the formula

$$A^\circ = \mathfrak{C}(\overline{\mathfrak{C}A}).$$

The operation $A \mapsto A^\circ$ is antiextensive, increasing, and idempotent. It is thus an opening in the algebraic sense, which of course means that it is a closure operator if we reverse the order: $A \leq B$ shall mean $A \supset B$. In addition to being an opening, it satisfies $X^\circ = X$ and $(A \cap B)^\circ = A^\circ \cap B^\circ$. Conversely, we may take these properties as axioms and define a set to be open if it is in the image of the operator. Then we get a topology and the operation of taking the interior of a set for this topology is equal to the original operator.

Summing up, we have five equivalent ways to define a topology: using open sets, closed sets, neighborhoods, taking the topological closure, and taking the interior.

If we have two topologies $\mathcal{U}_1(X)$ and $\mathcal{U}_2(X)$ on the same set X we say that the first is *weaker* or *coarser* than the second, and that the second is *finer* or *stronger* than the first, if $\mathcal{U}_1(X) \subset \mathcal{U}_2(X)$. Expressed in terms of closure operators, this means that $c_2 \leq c_1$ if c_j denotes the closure operator associated with $\mathcal{U}_j(X)$, $j = 1, 2$. The weakest topology is the *chaotic topology* $\{\emptyset, X\}$ and the strongest is the *discrete topology* $\mathcal{P}(X)$. The closure of a nonempty set in the chaotic topology is always the whole space, whereas the closure of a set in the discrete topology is the set itself.

A two-point space can have four topologies: in addition to the two just mentioned, they are $\{\emptyset, \{x\}, \{x, y\}\}$ and $\{\emptyset, \{y\}, \{x, y\}\}$. The two latter are called *Sierpiński topologies*.²⁰ Of the four, only three are different in the sense that they cannot be obtained from another one by renaming the points.

8.3. Transport of topologies

If $f: X \rightarrow Y$ is a mapping from a set X into a topological space Y we can transport the topology on Y to X by defining a subset of X to be open if and only if it is of the form $f^*(U)$ for some open subset U of Y . Because f^* is both a dilation and an erosion, it is clear that the family of all sets

$$(f^*)_*(\mathcal{U}(Y)) = \{f^*(U); U \in \mathcal{U}(Y)\}$$

²⁰Wacław Sierpiński, 1882–1969.

is a topology. Here we have used the notation introduced in (8.1) at the next higher level: $f^*: \mathcal{P}(Y) \rightarrow \mathcal{P}(X)$, $(f^*)_*: \mathcal{P}(\mathcal{P}(Y)) \rightarrow \mathcal{P}(\mathcal{P}(X))$; see (8.6). For brevity we shall denote $(f^*)_*(\tau) = (f^*)_*(\mathcal{U}(Y))$ by $f^-(\tau) = f^-(\mathcal{U}(Y))$, the *pull-back* of $\tau = \mathcal{U}(Y)$.

If $d: \mathcal{P}(Y) \rightarrow \mathcal{P}(Y)$ is a closure operator in Y , then $d^- = f^* \circ d \circ f_*$ is a closure operator in X , and if d satisfies the topological axioms $d(\emptyset) = \emptyset$ and $d(B_1 \cup B_2) = d(B_1) \cup d(B_2)$, then d^- does the same. Thus we can transport topological closure operators from Y to X . One can verify that the transported open sets correspond to the transported closure operator.

A particularly common case is when X is a subset of Y and f is the inclusion mapping. Then we say that the topology $f^-(\mathcal{U}(Y))$ defined on X is the *induced topology*. We see that U is open in X if and only if $U = V \cap Y$ for some open set in Y ; we also see that the closure operator d^- in X is defined as $d^-(A) = d(A) \cap X$.

If X is a topological space and $f: X \rightarrow Y$ a mapping of X into a set Y , we can of course consider the family $\{f_*(A); A \in \mathcal{U}(X)\} \subset \mathcal{P}(Y)$. However, this family is usually not a topology on Y , since f_* is not an erosion. Instead we use again f^* and declare a subset B of Y to be open if $f^*(B)$ is open in X . And we can verify that this is indeed a topology on Y ; we shall denote it by

$$f_-(\sigma) = f_-(\mathcal{U}(X)) = (f^*)^*(\mathcal{U}(X)) = \{B \in \mathcal{P}(Y); f^*(B) \in \mathcal{U}(X)\},$$

the *push-forward* of the topology $\sigma = \mathcal{U}(X)$ on X .

A common instance of this definition is when Y is a quotient set of X , i.e., when we have an equivalence relation \sim in X and let $Y = X/\sim$ be the set of all equivalence classes in X with respect to the relation. The mapping f associates to each element in X its equivalence class in Y . Then a subset B of Y is open in Y with respect to the topology we have pushed forward from X if and only if the union of all equivalence classes in B is open in X . The topology obtained in this way on X/\sim is called the *quotient topology*.

If $f: X \rightarrow Y$ is injective, and if we have a topology on X , push it forward to Y and then pull it back to X , the new topology agrees with the original one: $f^-(f_-(\mathcal{U}(X))) = \mathcal{U}(X)$. Similarly, if f is surjective and we start with a topology on Y , pull it back to X and then push it forward to Y , we obtain the original topology; $f_-(f^-(\mathcal{U}(Y))) = \mathcal{U}(Y)$. (This works so well because we did not use f_* but f^* in the definition of f_- .)

However, if we have a closure operator c in X , we cannot define a closure operator in Y by something like $c_- = f_* \circ c \circ f^*$.

8.4. Continuous mappings

Let $f: X \rightarrow Y$ be a mapping of a topological space X into a topological space Y and let x be a point in X . We say that f is *continuous at x* if $f^*(V)$ is a neighborhood of x for every neighborhood V of $f(x)$. It is called *continuous* if it is continuous at every point in X . We can translate this well-known notion into the language of open sets, closed sets, and closure operators. We can prove that f is continuous if and only if $f^*(U) \in \mathcal{U}(X)$ for every $U \in \mathcal{U}(Y)$; in other words if and only if the topology $f^-(\mathcal{U}(Y))$ is weaker than the topology $\mathcal{U}(X)$.

8.5. Connectedness

The family of all open and closed sets of a topological space X (sometimes called the *clopen sets*) is a Boolean algebra. This algebra must contain the two sets \emptyset, X , for they

are always both open and closed. (If X is empty, there is of course only one such set.) A topological space is said to be *connected* if it is nonempty and the only sets which are both open and closed are the empty set and the whole space. A subset of a topological space is called *connected* if it is connected as a topological space with the induced topology.²¹ A *connectivity component* (sometimes called a “connected component”) of a topological space is a connected subset which is maximal with respect to inclusion.

A connected subset which is both open and closed is a component. It is easy to prove that the closure of a connected subset is connected. Therefore all components are closed. They need not be open.

Proposition 8.1. *Let $f: X \rightarrow Y$ be a continuous mapping of a topological space X into a topological space Y . If X is connected, then so is $f_*(X) = \text{im} f$.*

Proof. Let B be a clopen subset of $\text{im} f$. Then $f^*(B)$ is clopen in X . Hence $f^*(B)$ is either empty or equal to X . Therefore $f_*(f^*(B)) = B \cap \text{im} f = B$ is either empty or equal to $\text{im} f$. This means that $\text{im} f$ is connected.

Corollary 8.2. *Let $f: X \rightarrow Y$ be a mapping of a topological space X into a set Y . Equip Y with the strongest topology such that f is continuous. Suppose that X is connected. Then $\text{im} f$ is connected, and the points in $Y \setminus \text{im} f$ are isolated. In particular, any quotient space of a connected topological space is connected.*

Proof. For any point $y \in Y$ not in the image of f , the inverse image $f^{-1}(y) = f^*({y})$ is empty, thus both open and closed. This means that $Y \setminus \text{im} f$ has the discrete topology and the connectivity components are just the singleton sets.

Of the four topologies that can live on a space consisting of two points, only three are connected, and out of these, only two are different in the sense that they cannot be obtained from another one by renaming.

8.6. Quotient spaces

In particular we shall use Corollary 8.2 with $X = \mathbf{R}$ and $Y = \mathbf{Z}$ to define connected topologies on the digital line \mathbf{Z} . Let $f: \mathbf{R} \rightarrow \mathbf{Z}$ be a surjective mapping. Then \mathbf{Z} equipped with the strongest topology such that f is continuous is a connected topological space. Thus we consider \mathbf{Z} as a quotient space of \mathbf{R} , not as a subspace. Now there exist very many surjective mappings $\mathbf{R} \rightarrow \mathbf{Z}$. It is not unnatural to restrict attention to increasing surjections $f: \mathbf{R} \rightarrow \mathbf{Z}$. Then $f^*({n})$ is an interval for every integer n ; denote its endpoints by a_n and $b_n \geq a_n$, so that

$$]a_n, b_n[\subset f^*({n}) \subset [a_n, b_n].$$

We can normalize the situation to $a_n = n - \frac{1}{2}$, $b_n = n + \frac{1}{2}$; this does not change the topology on \mathbf{Z} . Then $f(x) = \lfloor x + \frac{1}{2} \rfloor$ for all $x \in \mathbf{R} \setminus (\mathbf{Z} + \frac{1}{2})$, and $f(n + \frac{1}{2}) = n$ or $f(n + \frac{1}{2}) = n + 1$ for $n \in \mathbf{Z}$. The topology is therefore determined if we know for which n we have $f(n + \frac{1}{2}) = n$. For every subset A of \mathbf{Z} we get a topology on \mathbf{Z} by declaring that $f(n + \frac{1}{2}) = n$ for $n \in A$ and that, for all other real numbers x , we have

²¹According to Bourbaki (1961:I:§11:1) the empty space is connected. Here I follow instead the advice of Adrien Douady (personal communication, 2000-06-26), declaring the empty set not to be connected.

$f(x) = \lfloor x + \frac{1}{2} \rfloor$. Thus A describes faithfully all topologies obtained from increasing surjections—the others are just too many...

This can be explained as follows. It is natural to think of \mathbf{Z} as an approximation of the real line \mathbf{R} and to consider mappings $f: \mathbf{R} \rightarrow \mathbf{Z}$ expressing this idea. We may define $f(x)$ to be the integer closest to x ; this is well-defined unless x is a half-integer: $f(x) = \lfloor x + \frac{1}{2} \rfloor$ when $x \in \mathbf{R} \setminus (\mathbf{Z} + \frac{1}{2})$. So when $x = n + \frac{1}{2}$ we have a choice for each n : shall we define $f(n + \frac{1}{2}) = n$ or $f(n + \frac{1}{2}) = n + 1$? If we choose the first alternative for every n , thus putting $f^*({n}) =]n - \frac{1}{2}, n + \frac{1}{2}[$, the topology defined in Corollary 8.2 is called the *right topology* on \mathbf{Z} ; if we choose the second, we obtain the *left topology* on \mathbf{Z} ; cf. Bourbaki (1961:I:§1: Exerc. 2).

Another choice is to always choose an even integer as the best approximant of a half-integer. Then the closed interval $[-\frac{1}{2}, \frac{1}{2}]$ is mapped to 0, so $\{0\}$ is closed, whereas the inverse image of 1 is the open interval $] \frac{1}{2}, \frac{3}{2} [$, so that $\{1\}$ is open. This topology was introduced by E. D. Halimskiĭ (Efim Khalimsky), and we shall call it the *Khalimsky topology*; \mathbf{Z} with this topology is called the *Khalimsky line*. The Khalimsky line is connected, but the complement of any point is disconnected. Among all the topologies defined by increasing surjections $f: \mathbf{R} \rightarrow \mathbf{Z}$ only two have this property: the one just defined and the one obtained by translating everything by one step. For the left topology, for instance, all nonempty subsets are connected.

8.7. Separation axioms

The closure of a subset A of a topological space X will be denoted by \overline{A} . The intersection of all neighborhoods of a point x will be denoted by $N(x)$. We note that $x \in \overline{\{y\}}$ if and only if $y \in N(x)$. The relation $x \in \overline{\{y\}}$ defines a preorder in X , i.e., a relation satisfying (3.1) and (3.2) but not necessarily (3.3). We shall denote it by $x \preceq y$; thus $x \preceq y$ if and only if $x \in \overline{\{y\}}$. It was introduced by Aleksandrov (1937:503). We shall call it the *specialization preorder*; cf. Kong et al. (1991:905).

A *Kolmogorov space* (Bourbaki (1961:I:§1: Exerc. 2)), also called a T_0 -space, is a topological space such that $x \in N(y)$ and $y \in N(x)$ only if $x = y$, thus precisely when the specialization preorder is an order (satisfies (3.3)). It is quite reasonable to impose this axiom; if x belongs to the closure of $\{y\}$ and vice versa, then x and y are indistinguishable from the point of view of topology, we cannot distinguish points from knowledge of the open sets to which they belong. We should therefore identify them and consider a quotient space.

So every topology τ on a set X defines a preorder $P(\tau) = \preceq$ in X , and this preorder is an order if and only if τ satisfies Kolmogorov's axiom.

Conversely, every preorder \preceq in a set X defines a topology $T(\preceq)$ on X whose smallest neighborhoods are

$$N(x) = \{y \in X; x \preceq y\}.$$

This means that a subset A of X is open if and only if, for all $a \in A$ and all $y \in X$, $a \preceq y$ implies $y \in A$.

If we start with a preorder \preceq , define the topology $T(\preceq)$, and then define the preorder $P(T(\preceq))$, then we get back to the original preorder: $\preceq = P(T(\preceq))$; in other words $P \circ T$ is the identity. However, if we start with a topology τ , construct the preorder $P(\tau)$, then the topology $T(P(\tau))$ defined by that preorder is not equal to τ in general; $T \circ P$ is

not the identity. However, $T(P(\tau)) = \tau$ if and only if τ defines a smallest-neighborhood space as defined in the next section.

The separation axiom T_1 states that $N(x) = \{x\}$. It is too strong to be of interest for the spaces considered here. The specialization preorder in this case is the discrete order: we have $x \preceq y$ if and only if $x = y$.

Two points x and y in a topological space Y are said to be *adjacent* if $x \neq y$ and $\{x, y\}$ is connected. We note that $\{x, y\}$ is connected if and only if either $x \in N(y)$ or $y \in N(x)$.

8.8. Smallest neighborhood spaces

In a topological space the union of any family of open sets is open. It may happen that also the intersection of any family of open sets is open. Equivalently, every point in the space possesses a smallest neighborhood. A topological space with this property we shall call a *smallest-neighborhood space*. Another suitable name would be a *P. S. Aleksandrov space*, in honor of P. S. Aleksandrov,²² who introduced them in (1935, 1937). It is equivalent to require that the union of an arbitrary family of closed sets be closed.

The intersection $N(x)$ of all neighborhoods of a point x is open for all x if and only if the space is a smallest-neighborhood space.

Aleksandrov (1935, 1937) introduced the term *espace discret, diskreter Raum (discrete space)* for a topological space such that the intersection of any family of open sets is open. The intersection of all closed sets containing a set M he called its *Hülle (hull)*, and denoted it by \overline{M} or AM . The intersection of all open sets containing a set M he called its *Stern (star)* and denoted it by OM . He noted that the star of a set is a closure operation satisfying the two extra conditions $\overline{\emptyset} = \emptyset$ and $\overline{A \cup B} = \overline{A} \cup \overline{B}$ (see the discussion in section 3), and therefore defines a topology, which he called *reciproque* (1935) or *dual* (1937). The closed set of a smallest-neighborhood space satisfies the axioms of the open sets of a topology, so there is a complete symmetry between the two topologies in such a space.

It is easy to see that a mapping $f: X \rightarrow Y$ between two smallest-neighborhood spaces is continuous if and only if it is increasing for the specialization preorder. Thus continuity in these spaces is actually order theoretic, and the smallest-neighborhood spaces are actually special cases of preordered sets. This means that the rich theory of (pre)ordered sets can be put to work here.

Aleksandrov's choice of terms seems fortunate, but nowadays it is not possible to use the term *discrete space* in Aleksandrov's sense, since the discrete topology in modern usage refers only to the topology where every set is open, the strongest of all topologies. This is why I propose to call a discrete space in Aleksandrov's sense a *smallest-neighborhood space* or a *P. S. Aleksandrov space*.

The closed points, i.e., the points x such that $\overline{\{x\}} = \{x\}$, Aleksandrov called *Eckpunkte (vertices)*, and the open points, i.e., the points x such that the singleton $\{x\}$ is open, he called *Grundpunkte (base points)*.

We can define a topology on the digital line \mathbf{Z} by declaring all odd points to be open, thus $N(2k+1) = \{2k+1\}$, and all even points to have a smallest neighborhood $N(2k) = \{2k-1, 2k, 2k+1\}$. It follows that the even points are closed, for the

²²Pavel Sergeevič Aleksandrov, 1896–1982; not to be confused with Aleksandr Danilovič Aleksandrov, 1912–1999.

complement of an even point $2k$ is the union of all $N(x)$ with $x \neq 2k$, thus an open set. This is the Khalimsky topology already defined in Chapter 8.5. Thus in the Khalimsky topology the even points are *Eckpunkte* and the odd points are *Grundpunkte* in Aleksandrov's terminology. In the specialization order, the base points are higher than the vertices....??

A *Khalimsky interval* is an interval $[a, b] \cap \mathbf{Z}$ equipped with the topology induced by the Khalimsky topology on \mathbf{Z} . A *Khalimsky circle* is a quotient space $\mathbf{Z}_m = \mathbf{Z}/m\mathbf{Z}$ of the Khalimsky line for some even integer $m \geq 4$. (If m is odd, the quotient space receives the chaotic topology, which is not interesting.)

The *Khalimsky plane* is the Cartesian product of two Khalimsky lines, and more generally, *Khalimsky n -space* is the Cartesian product of n copies of \mathbf{Z} . Equivalently, we can define Khalimsky n -space by declaring $\{x \in \mathbf{Z}^n; \|x - c\|_\infty \leq 1\}$ to be open for any point $c \in (2\mathbf{Z})^n$ and then taking all intersections of such sets as open sets, then all unions of such intersections.

There are, however, other topologies in \mathbf{Z}^2 which are of interest: we may declare $\{x \in \mathbf{Z}^2; \|x - c\|_1 \leq 1\}$ to be open for any c such that $\sum c_j \in 2\mathbf{Z}$ as well as all intersections of such sets.²³ The Khalimsky topology and the topology just defined are not comparable: none is stronger than the other. However, they are related, for if we turn the Khalimsky plane 45° and delete all points which are not open or closed, we obtain the new topology.

To exhibit some of the analogies between topological spaces and preordered sets, let us list some properties of continuous and increasing mappings.

| <i>Continuous mappings $X \rightarrow Y$ between topological spaces</i> | <i>Increasing mappings $X \rightarrow Y$ between preordered sets</i> |
|--|---|
| X has the discrete topology \Rightarrow all mappings are continuous | X has the discrete order \Rightarrow all mappings are increasing |
| Y has the chaotic topology \Rightarrow all mappings are continuous | Y has the chaotic preorder \Rightarrow all mappings are increasing |
| X has the chaotic topology and Y has a Kolmogorov topology \Rightarrow only the constants are continuous | X has the chaotic preorder and Y is ordered \Rightarrow only the constants are increasing |
| Y has the discrete topology and X is connected \Rightarrow only the constants are continuous | Y has the discrete order and X is connected \Rightarrow only the constants are increasing |

8.9. Exercises

8.1. Prove that if $f: X \rightarrow Y$ is a mapping of a set X into a set Y , then the closing $f^*/_*f^*$ is equal to $f^* \circ f_* \geq \text{Id}_{\mathcal{P}(X)}$. Similarly $f_*/_*f_* = f_* \circ f^*$ (an opening; cf. Theorem 7.18).

8.2. Show that if we have a closure operator c in X and a mapping $f: X \rightarrow Y$, then we cannot define a closure operator in Y by something like $c_- = f_* \circ c \circ f^*$. In general c_- will be

²³I found this topology in response to a question asked by Timur Sadykov on February 1, 2000. However, I found out later that it was defined already by Wyse (1970).

neither extensive nor idempotent. Construct examples which show this. How shall we define the closure operator connected with the topology $f_-(\mathcal{U}(X))$ on Y ?

8.3. Express the continuity of a mapping $f: X \rightarrow Y$ from one topological space X into another, Y , in terms of the families of all closed sets $\mathcal{F}(X)$ and $\mathcal{F}(Y)$ as well as in terms of the closure operators in X and Y .

8.4. We have seen that a two-point space can have four topologies. How many topologies are there on a three-point space? How many of these are different in the sense that they cannot be obtained from another one by renaming the points? How many are connected? How many of these are different?

9. A closer look at the Khalimsky plane

9.1. Continuous functions

In this chapter we shall take a closer look at the Khalimsky plane, and in particular consider the Jordan curve theorem and Brouwer's fixed-point theorem in the new setting. We first explain the meaning of continuity.

We recall that a subset A of \mathbf{Z} is open for the Khalimsky topology if and only if, for every even number $2n \in A$, also the two odd numbers $2n \pm 1$ belong to A . To construct the Khalimsky plane, we take the Cartesian product of two copies of the Khalimsky line. The topology is then determined by the rule that a subset A of the Khalimsky plane is open if and only if, for every pair of even numbers $x = (2m, 2n)$, all pairs $y \in \mathbf{Z}^2$ with $\|y - x\|_\infty \leq 1$ belong to A , for every pair $(2m, 2n+1) \in A$ also $(2m \pm 1, 2n+1) \in A$, and, finally, for every pair $(2m+1, 2n) \in A$, also $(2m+1, 2n \pm 1) \in A$.

A function $f: X \rightarrow Y$ from one smallest-neighborhood space into another is continuous at a point x if and only if the direct image of $N_X(x)$ is contained in $N_Y(f(x))$, or, equivalently, the inverse image of $N_Y(f(x))$ contains $N_X(x)$:

$$(9.1) \quad f_*(N_X(x)) \subset N_Y(f(x)), \text{ equivalently } N_X(x) \subset f^*(N_Y(f(x))).$$

Here $N_X(x)$ and $N_Y(y)$ denote the smallest neighborhood of $x \in X$ and $y \in Y$, respectively. If we apply this to the case when $X = Y = \mathbf{Z}$, it means the following. If $x \in \mathbf{Z}$ is odd, the property always holds; if x is even and $f(x)$ is odd, it means that $f(x \pm 1) = f(x)$, and if x is even and also $f(x)$ is even, it means that $|f(x \pm 1) - f(x)| \leq 1$. In particular, a continuous function is Lip-1, but it must sometimes have intervals of constancy, viz. every time it takes an odd value at an even point (and hence also when it takes an even value at an odd point).

We observe that the following functions are continuous: (1) $\mathbf{Z} \ni x \mapsto a \in \mathbf{Z}$, where a is a constant; (2) $\mathbf{Z} \ni x \mapsto \pm x + c \in \mathbf{Z}$, where c is an even constant; (3) $\max(f, g)$ and $\min(f, g)$ if f, g are continuous. Actually every continuous function on a bounded Khalimsky interval can be obtained by a finite succession of the rules (1), (2), (3). Note that the function $x \mapsto x + 1$ is discontinuous.

For functions of two variables, $f: \mathbf{Z}^2 \rightarrow \mathbf{Z}$, (9.1) means the following. We list the eight possible parities of the triple $(x_1, x_2, f(x)) \in \mathbf{Z}^3$:

| Parity of x_1 | Parity of x_2 | Parity of $f(x)$ | Condition |
|-----------------|-----------------|------------------|--|
| Odd | Odd | Odd | None |
| Odd | Odd | Even | None |
| Odd | Even | Odd | $y_1 = x_1, y_2 - x_2 \leq 1 \Rightarrow f(y) = f(x)$ |
| Odd | Even | Even | $y_1 = x_1, y_2 - x_2 \leq 1 \Rightarrow f(y) - f(x) \leq 1$ |
| Even | Odd | Odd | $ y_1 - x_1 \leq 1, y_2 = x_2 \Rightarrow f(y) = f(x)$ |
| Even | Odd | Even | $ y_1 - x_1 \leq 1, y_2 = x_2 \Rightarrow f(y) - f(x) \leq 1$ |
| Even | Even | Odd | $\ y - x\ _\infty \leq 1 \Rightarrow f(y) = f(x)$ |
| Even | Even | Even | $\ y - x\ _\infty \leq 1 \Rightarrow f(y) - f(x) \leq 1$ |

This table lists what it means for f to be continuous at a particular point x . However, the description becomes much simpler if we use the specialization order \preceq , for then a continuous function is just an increasing function. We know that, in \mathbf{Z} ,

$$\dots \preceq -2 \succ -1 \preceq 0 \succ 1 \preceq 2 \succ 3 \preceq 4 \succ \dots$$

In \mathbf{Z}^2 , $(0, 0) \succ (1, 0), (0, 1) \succ (1, 1)$ and, in general,

$$(2m, 2n) \succ (2m + 1, 2n), (2m, 2n + 1) \succ (2m + 1, 2n + 1) \text{ for all } m, n \in \mathbf{Z}.$$

So continuity at x boils down to $x \preceq y \Rightarrow f(x) \preceq f(y)$ for all $y \in \mathbf{Z}^2$; continuity everywhere to the same implication but now for all $x, y \in \mathbf{Z}^2$. For example, if both components of x are odd, the only y which satisfies $x \preceq y$ is $y = x$, so $f(x) \preceq f(y)$ holds automatically. If, on the other hand both components of x are even, then $\{y; x \preceq y\} = B_{\preceq}(x, 1)$ for the l^∞ norm, and if $f(x)$ in addition is odd, then $f(x) \preceq f(y)$ holds only for $f(y) = f(x)$, so f must be constant on $B_{\preceq}(x, 1)$.

We note that if $x, y \in \mathbf{Z}$ and $x \preceq y$, then $|x - y| \leq 1$. Conversely, if $|x - y| \leq 1$, then either $x \preceq y$ or $y \preceq x$. Hence $|x - y| \leq 1$ implies $|f(x) - f(y)| \leq 1$ for any continuous function $f: \mathbf{Z} \rightarrow \mathbf{Z}$, and we see that f is Lip-1. In two variables we have the same conclusion. In the proof of this fact we shall need the following notation. For any two points $x, y \in \mathbf{Z}^2$ we define $q(x, y) = (x_1, y_2)$. The four points $x, y, q(x, y), q(y, x)$ thus form a rectangle (perhaps degenerate); if $y_j = x_j \pm 1, j = 1, 2$, they form a square with side length 1.

Theorem 9.1. *A continuous function $f: \mathbf{Z}^2 \rightarrow \mathbf{Z}$ is Lip-1 for the l^∞ norm. More generally, the conclusion holds for any continuous function $f: X \rightarrow \mathbf{Z}$, where X is a connected subset of \mathbf{Z}^2 such that $q(x, y), q(y, x) \in X$ for all $x, y \in X$ such that $y_j = x_j \pm 1, j = 1, 2$, and we do not have $x \preceq y$, nor $y \preceq x$.*

Proof. It is enough to prove that if $\|x - y\|_\infty \leq 1$, then $|f(x) - f(y)| \leq 1$. If $x, y \in X$ with $x \preceq y$, then $f(x) \preceq f(y)$, which in turn implies that $|f(x) - f(y)| \leq 1$. Assume now that $\|x - y\|_\infty \leq 1$ but that we do not have $x \preceq y$, nor $y \preceq x$. Then we have essentially the case $x = (1, 0), y = (0, 1)$. We have $q(y, x) = (0, 0) \succ x, y \succ (1, 1) = q(x, y)$, which implies that $f(0, 0) \succ f(x), f(y) \succ f(1, 1)$. (By hypothesis both $(0, 0)$ and $(1, 1)$ belong to X .) Hence $|f(0, 0) - f(1, 1)| \leq 1$. Also, $|f(0, 0) - f(x)| \leq 1$ and $|f(x) - f(1, 1)| \leq 1$, which implies that $f(x) = f(0, 0)$ or $f(1, 1)$. Similarly, $f(y) = f(0, 0)$ or $f(y) = f(1, 1)$. Therefore $|f(x) - f(y)| \leq 1$.

This result holds for many subsets X of \mathbf{Z}^2 , but not for all:

Example. Let $X = \{x \in \mathbf{Z}^2; \|x\|_1 \leq 1\}$, the l^1 ball in \mathbf{Z}^2 and define $f(x) = x_1 - x_2$. Then X is connected for the Khalimsky topology and f is continuous, but f is not Lip-1. We note that $q(y, x) = (0, 0) \in X$ but that $q(x, y) = (1, 1) \notin X$ if we take $x = (1, 0) \in X$, $y = (0, 1) \in X$. So both $q(x, y)$ and $q(y, x)$ need to be elements of X .

Theorem 9.2. *A function $f: \mathbf{Z}^2 \rightarrow \mathbf{Z}$ is continuous if and only if it is separately continuous. More generally, the equivalence holds for any function $f: X \rightarrow \mathbf{Z}$ where X is a subset of \mathbf{Z}^2 such that one of $q(x, y)$, $q(y, x)$ belongs to X if $y_j = x_j \pm 1$ and $x \preceq y$.*

Proof. Assume that f is separately continuous and that $x \preceq y$. Then we shall prove that $f(x) \preceq f(y)$. If $x_1 = y_1$, then $x_2 \preceq y_2$, and the inequality $f(x) \preceq f(y)$ follows from the separate continuity of the function $x_2 \mapsto f(x)$ for a fixed x_1 . The conclusion is similar if $x_2 = y_2$; then the continuity of $x_1 \mapsto f(x)$ for a fixed x_2 does the job. The case when $x_1 \neq y_1$ and $x_2 \neq y_2$ remains to be considered. Then $y_j = x_j \pm 1$. One of the points $q(x, y)$ and $q(y, x)$ belongs to X ; let z be one of them that does. Then clearly $x \preceq z \preceq y$, which in view of the separate continuity implies $f(x) \preceq f(z)$ and $f(z) \preceq f(y)$, and we are done.

Example. Let $X = \{0\} \cup \{x \in \mathbf{Z}^2; |x_1| = |x_2| = 1\}$. This set is connected for the Khalimsky topology. Every function $f: X \rightarrow \mathbf{Z}$ is separately continuous, but not all of them are continuous. With $x = (0, 0)$ and $y = (1, 1)$ we see that none of $q(x, y) = (0, 1)$ and $q(y, x) = (1, 0)$ belongs to X .

Like in real analysis there is an intermediate-value theorem for the Khalimsky line:

Theorem 9.3. *Let two continuous functions $f, g: I \rightarrow \mathbf{Z}$ be given on a Khalimsky interval $I = [a, b] \cap \mathbf{Z}$. Assume that there are points $s, t \in I$ with $f(s) \geq g(s)$ and $f(t) \leq g(t)$. Then there exists a point p , intermediate between s and t , such that $f(p) = g(p)$.*

Proof. Without loss of generality we may assume that $s \leq t$. Define

$$M = \{x \in \mathbf{Z}; s \leq x \leq t \text{ and } f(x) \geq g(x)\}.$$

Clearly $s \in M$, so M is not empty. Let $p = \max_{x \in M} x$. If $p = t$, then $f(t) = g(t)$ and we are done; if not, $p + 1 \leq t$. Then we must have $f(p) \geq g(p)$, $f(p + 1) < g(p + 1)$. We claim that $f(p) = g(p)$. If this were not true, we would have $f(p) \geq g(p) + 1$, $f(p + 1) \leq g(p + 1) - 1$. Because of the Lipschitz continuity, the only possibility then would be $f(p) = g(p) + 1$ and $f(p + 1) = g(p + 1) - 1$. But even this situation is impossible. If p and $f(p)$ are of different parity, then $f(p + 1) = f(p)$, which would require a jump of two units in g . If on the other hand p and $f(p)$ are of the same parity, then p and $g(p)$ are of different parity, so that $g(p + 1) = g(p)$, requiring a jump of two units in f . This contradiction leaves us with the only possibility $f(p) = g(p)$.

9.2. A fixed-point theorem

Let us say that a topological space X has the *fixed-point property* if every continuous mapping $f: X \rightarrow X$ possesses a *fixed point*, i.e., a point p such that $f(p) = p$. The

famous Brouwer²⁴ fixed-point theorem states that a closed ball in \mathbf{R}^n has the fixed-point property. This theorem is a deep result when $n \geq 2$ but easy to prove if $n = 1$: we may take

$$(9.2) \quad p = \sup_x (x; f(x) \geq x).$$

If $f: [0, 1] \rightarrow [0, 1]$ is a mapping from a compact interval into itself which is increasing for the usual order inherited from the real axis, then it also has a fixed point, and formula (9.2) again gives a fixed point. This generalizes to complete lattices; if $f: L \rightarrow L$ is an increasing mapping of a complete lattice L into itself, then the point p defined by (9.2) is a fixed point of f . This is Tarski's fixed point theorem (1955), exhibiting one of the many analogies between continuous mappings and increasing mappings. However, a Khalimsky interval of more than two points is not a lattice.

In this section we shall prove that certain subsets of the Khalimsky plane possess the fixed-point property. Since the Khalimsky line and the Khalimsky plane are ordered sets and the continuous mappings are precisely the increasing mappings, there are theorems from the theory of ordered sets that can be applied here; cf. Baclawski & Björner (1979). However, the proofs presented here are more direct in the context of digital geometry.

A topological space cannot have the fixed-point property unless it is connected. Indeed, if U and V are two disjoint nonempty open sets whose union is the whole space X , then we may define $f(x) = b$ for all $x \in U$ and $f(x) = a$ for all $x \in V$, where a is an arbitrary point in U and $b \in V$. Then all inverse images are open, for they are either empty, equal to U , equal to V , or equal to X . Thus f is certainly continuous, but it has no fixed point.

On the other hand, connectedness is by no means sufficient. The mapping $f(x) = -x$ is continuous in \mathbf{R}^n , and in fact in any abelian group with a reasonable topology; it has no fixed point in X if $0 \notin X$. Any nonempty subset X of such a group with $0 \notin X$ and such that $x \in X$ implies $-x \in X$ will be a space such that some continuous selfmapping is without fixed point.

In a finite set X with N points there are N^N selfmappings $X \rightarrow X$. Out of these, $(N - 1)^N$ do not have fixed points; there are thus $N^N - (N - 1)^N$ mappings which have a fixed point. The proportion of mappings with a fixed point is $1 - (1 - 1/N)^N$, which tends to $1 - 1/e$ (about 63 percent) when N tends to infinity. If we introduce a topology on X we may ask how many of the N^N mappings are continuous. Let us denote by C the number of continuous mappings. If the space possesses the fixed-point property the number of mappings of different kinds can then be listed as follows. The table contains only one unknown, C , the number of continuous mappings.

| | <i>Continuous</i> | <i>Discontinuous</i> | <i>Sum</i> |
|-----------------------|-------------------|-----------------------|-------------------|
| <i>Fixed point</i> | C | $N^N - (N - 1)^N - C$ | $N^N - (N - 1)^N$ |
| <i>No fixed point</i> | 0 | $(N - 1)^N$ | $(N - 1)^N$ |
| <i>Sum</i> | C | $N^N - C$ | N^N |

²⁴Luitzen Egbertus Jan Brouwer, 1881–1966.

We shall now prove that a continuous mapping of an interval into itself has a fixed point. Later we shall consider subsets of the Khalimsky plane, which we shall attack using an implicit-function theorem.

Theorem 9.4. *Let us define a subset A of the set of all continuous mappings of \mathbf{Z} into itself,*

$$(9.3) \quad A = \{f \in C(\mathbf{Z}, \mathbf{Z}); \exists s \in \mathbf{Z}, f(s) \geq s \text{ and } \exists t \in \mathbf{Z}, f(t) \leq t\}.$$

Then $f \in C(\mathbf{Z}, \mathbf{Z})$ has a fixed point if and only if $f \in A$.

Proof. If f has a fixed point, it is obvious that f belongs to A . The converse is a special case of the intermediate-value theorem (Theorem 9.3), taking $g(x) = x$.

Corollary 9.5. *Every bounded Khalimsky interval has the fixed-point property.*

Proof. Let $f: I \rightarrow I$ be a continuous mapping, where $I = [a, b] \cap \mathbf{Z}$ is a bounded interval. Extend f to a mapping $g: \mathbf{Z} \rightarrow \mathbf{Z}$ by defining $g(x) = f(a)$ for $x < a$ and $g(x) = f(b)$ for $x > b$. Then it is easy to check that g is continuous, and it is obvious that g belongs to A . Thus it has a fixed point $p \in \mathbf{Z}$, but as $p \in \text{img} \subset I$, p is a fixed point also of f .

Example. For a Khalimsky interval consisting of two points only, say $\{a, a + 1\}$, there are four mappings: the two constant mappings, the identity, and the one interchanging a and $a + 1$. The first three obviously have a fixed point; the fourth does not. But it is discontinuous. Thus the statistics looks like this ($N = 2, C = 3$):

| | <i>Continuous</i> | <i>Discontinuous</i> | <i>Sum</i> |
|-----------------------|-------------------|----------------------|------------|
| <i>Fixed point</i> | 3 | 0 | 3 |
| <i>No fixed point</i> | 0 | 1 | 1 |
| <i>Sum</i> | 3 | 1 | 4 |

Theorem 9.6. *Every bounded Khalimsky rectangle $I \times J \subset \mathbf{Z}^2$ has the fixed-point property.*

We shall prove this result using an implicit-function theorem. In real analysis there is an implicit-function theorem which says the following. If f is a real-valued function of class C^1 defined in an open subset Ω of \mathbf{R}^2 with a zero $a \in \Omega$, and its derivative with respect to the first variable, $\partial f / \partial x_1$, is non-zero at a , then there is a function g of one variable defined near a_2 such that $f(g(x_2), x_2) = 0$ for all x_2 near a_2 . This function g is also of class C^1 . One says that the equation $f(x_1, x_2) = 0$ defines x_1 *implicitly* as a function of x_2 ; the equation $f(x_1, x_2) = 0$ is equivalent to the *explicit* formula $x_1 = g(x_2)$ near a . The result is local in the sense that we only assert something in a neighborhood of a .

In the digital plane we have an implicit-function theorem of a similar kind. We shall solve the implicit equation $f(x, y) = x$. (In real analysis, this is equivalent to $f(x, y) - x = 0$, but the left-hand side $f(x, y) - x$ is not continuous in general in the digital case even if f is.)

Theorem 9.7. *Let us define an order, and hence a topology, in the space A of mappings defined by (9.3), by declaring that $f \preceq g$ if and only if $f(x) \preceq g(x)$ for all $x \in \mathbf{Z}$. Let us also define a mapping $F: A \rightarrow A$ by*

$$(9.4) \quad F(g) = \max(x \in \mathbf{Z}; g(x) \geq x), \quad g \in A;$$

$F(g)$ is thus the largest fixed point of g . The mapping F is continuous.

Lemma 9.8. *The mapping F defined by (9.4) satisfies*

$$(9.5) \quad g(x) \geq x \text{ if and only if } F(g) \geq x, \quad x \in \mathbf{Z}, g \in A.$$

Proof. The implication to the right follows easily from the definition of F . For the other implication we use the Lip-1 property of g : if $g(x) < x$, then $g(x') < x'$ for all $x' \geq x$, which implies that $F(g) < x$.

Proof of Theorem 9.7. We shall prove that $f \preceq g$ implies $F(f) \preceq F(g)$. Fix two functions f and g satisfying $f \preceq g$. We consider the two cases $F(f)$ odd, $F(f)$ even.

Case 1. Assume that $F(f)$ is odd. Without loss of generality we may assume that $F(f) = 1$, thus that $f(1) = 1$ and that $f(2) = 1$ or $f(2) = 0$. We shall prove that $F(g)$ is 0, 1, or 2. We write the inequality $0 \leq F(g) \leq 2$ as $0 \leq F(g) < 3$ and then translate it using Lemma 9.8 as $0 \leq g(0)$ and $g(3) < 3$. We know that $g(1) \succeq f(1) = 1$, which implies that $g(1) \geq 0$ and that $g(0) \geq 0$. Hence it only remains to be proved that $g(3) < 3$.

Case 1.1. Assume that $f(2) = 1$. Then $f(1) = f(2) = f(3) = 1$, so that $g(3) \succeq f(3) = 1$ and $g(3) \leq 2 < 3$.

Case 1.2. Assume that $f(2) = 0$. Then $g(2) \succeq f(2) = 0$, which implies that $g(2) = 0$; by Lipschitz continuity $g(3) \leq 1 < 3$.

Case 2. If $F(f)$ is even, we may assume that $F(f) = 0$ without loss of generality. We shall then prove that $F(g) = 0$. We write this as $0 \leq F(g) < 1$ and translate it using Lemma 9.8 as $0 \leq g(0)$ and $g(1) < 1$. From $F(f) = 0$ it follows that $f(0) = 0$ and that $f(1) = 0$ or $f(1) = -1$. Since $0 = f(0) \preceq g(0)$, it follows that $g(0) = 0$. It remains to be proved that $g(1) < 1$.

Case 2.1. Assume that $f(1) = -1$. Then $g(1) \succeq f(1) = -1$, which implies that $g(1) \leq 0 < 1$.

Case 2.2. Assume finally that $f(1) = 0$. Then $g(1) \succeq f(1) = 0$, so $g(1) = 0 < 1$.

Corollary 9.9. *Let $f: I \times J \rightarrow I$ be a continuous function defined in a rectangle $I \times J$, where I and J are Khalimsky intervals, I being bounded. Then*

$$(9.6) \quad h(y) = \max_{x \in I}(x; f(x, y) \geq x), \quad y \in J,$$

is continuous in J . The function h satisfies $f(h(y), y) = h(y)$.

Proof. The function h is given by $h(y) = F(x \mapsto f(x, y))$. It depends continuously on the function $x \mapsto f(x, y)$, which in turn is a continuous function of y .

Example. It is natural to ask if the function

$$h(y) = \max_{x \in I}(x; f(x, y) \geq g(x, y)), \quad y \in J,$$

is continuous for more general choices of functions f and g . That this is not always so is shown by the example $f(x, y) = y$, $g(x, y) = \min(x, 0)$. Then $h(y) = y$ if $y < 0$ and $h(y) = b$ if $y \geq 0$, assuming the interval I to be $[a, b] \cap \mathbf{Z}$. This yields a discontinuous function if $a \leq -1$, $b \geq 1$.

Proof of Theorem 9.6 The mapping f has two components f_1 and f_2 . As in (9.6) we let $h(y)$ denote the largest fixed point of the partial mapping $x \mapsto f_1(x, y)$ for a fixed y , $h(y) = F(x \mapsto f_1(x, y))$. So $f_1(h(y), y) = h(y)$ for all $y \in J$. We then form the composition $k(y) = f_2(h(y), y)$. As a composition of continuous mappings it is continuous, and by the one-dimensional theorem it possesses a fixed point q , $k(q) = q$. Collecting what we have, we see that $f_1(h(q), q) = h(q)$ and that $k(q) = f_2(h(q), q) = q$, which means that we have proved that $f(h(q), q) = (f_1(h(q), q), f_2(h(q), q)) = (h(q), q)$, thus that $(h(q), q)$ is a fixed point.

Example. The Khalimsky square $\{0, 1\}^2 \subset \mathbf{Z}^2$. There are $N^N = 4^4 = 256$ selfmappings of $\{0, 1\}^2$, of which $(N - 1)^N = 3^4 = 81$ do not have fixed points. The remaining $256 - 81 = 175$ have a fixed point. Of the 16 mappings $\{0, 1\}^2 \rightarrow \{0, 1\}$, 6 are continuous. There are therefore $6^2 = 36$ continuous mappings $\{0, 1\}^2 \rightarrow \{0, 1\}^2$, and we know already that they all have fixed points.

The table of different kinds of mappings therefore looks as follows ($N = 4$, $C = 6^2$).

| | <i>Continuous</i> | <i>Discontinuous</i> | <i>Sum</i> |
|-----------------------|-------------------|----------------------|------------|
| <i>Fixed point</i> | 36 | 139 | 175 |
| <i>No fixed point</i> | 0 | 81 | 81 |
| <i>Sum</i> | 36 | 220 | 256 |

In this simple case we can of course see directly that all continuous mappings have a fixed point. Indeed, of the 6 continuous mappings $\{0, 1\}^2 \rightarrow \{0, 1\}$, five map $(0, 0)$ to 0; the remaining one is the constant 1. Therefore, of the 36 continuous mappings $\{0, 1\}^2 \rightarrow \{0, 1\}^2$, $(0, 0)$ is a fixed point except when one of the components is the constant 1. Thus they all have a fixed point.

It is easy to generalize the implicit-function theorem and the fixed-point theorem to somewhat more general sets. We formulate an example of the latter:

Theorem 9.10. *Let X be a subset of \mathbf{Z}^2 defined as*

$$X = \{(x, y) \in \mathbf{Z} \times J; \varphi(y) \leq x \leq \psi(y)\},$$

where J is a bounded Khalimsky interval, and φ and ψ two continuous functions defined on J . Assume that $\varphi(y) < \psi(y)$, $y \in J$. Then X has the fixed-point property.

Proof. Take an interval $I = [a, b] \cap \mathbf{Z}$ which is so large that $\varphi(y), \psi(y) \in I$ for all $y \in J$. Then extend $f: X \rightarrow X$ to a mapping $g: I \times J \rightarrow I \times J$ by defining $g(x, y) = \varphi(y)$ when $a \leq x < \varphi(y)$ and $g(x, y) = \psi(y)$ when $\psi(y) < x \leq b$. Then g is continuous. It must have a fixed point in $I \times J$ according to Theorem 9.5. However, the fixed point must actually lie in X and be a fixed point of f .

The result on separate continuity (Theorem 9.2) makes it easy to go up in dimension.

9.3. Jordan curve theorems

There is a Jordan curve theorem in the Khalimsky plane. During the lectures I discussed this; here [at least in this version of the notes] I just refer to Kiselman (2000) and the references therein.

9.4. Exercises

9.1. How many continuous selfmappings are there in a Khalimsky interval with three elements?

9.2. How many continuous selfmappings are there in a Khalimsky interval with four elements like $\{0, 1, 2, 3\}$?

9.3. Prove in detail that there are six continuous mappings $\{0, 1\}^2 \rightarrow \{0, 1\}$ and hence that there are 36 continuous mappings $\{0, 1\}^2 \rightarrow \{0, 1\}^2$. Compare with the previous exercise; both spaces have four elements.

9.4. Try to estimate the number C of continuous mappings $I \rightarrow J$ between Khalimsky intervals, and more generally mappings $I_1 \times \cdots \times I_m \rightarrow J_1 \times \cdots \times J_m$ between boxes in Khalimsky spaces. Compare with the number of Lip-1 mappings.

10. Digitization

Digital geometry is about geometry in digital spaces—however, we shall not at this point give a formal definition. Suffice it to say that *digital* is used here as opposed to *continuous*. The space \mathbf{R}^n is a space where we do continuous geometry; the space \mathbf{Z}^n is an example of a space where we do digital geometry.²⁵

In this section we shall first discuss what a good digitization should mean, and then study the notion of a digital line.

10.1. What is digitization?

Let X be a set and Z an arbitrary subset of X . (Think of X as \mathbf{R}^n and Z as \mathbf{Z}^n if you want a basic example.) If we want to digitize X we may start with a mapping $f: X \rightarrow Z$ and then define the digitization of a set A as $f_*(A) = \{f(x); x \in A\}$; cf. (8.1).

However, this approach is too narrow; it is often not possible to start with a pointwise mapping. Instead, we shall define here a digitization of X into Z as a mapping $F: \mathcal{P}(X) \rightarrow \mathcal{P}(Z)$ with certain desirable properties. We shall think of $F(A)$ as a digital representation of A . A very simple such representation is $F(A) = A \cap Z$, but it is not very faithful, since many sets are mapped to the empty set, for example $A = X \setminus Z$. (However, it might work for sufficiently fat sets.) One desirable condition is therefore that $F(A)$ be empty only if A is empty. We also remark that this mapping is not of the form $F = f_*$.

We recall that the mappings f_* are dilations; see (8.5). It seems desirable to require in general that a digitization $F: \mathcal{P}(X) \rightarrow \mathcal{P}(Z)$ shall be a dilation. In particular this

²⁵The word itself comes the Latin *digitus*, meaning ‘finger, toe’ and being related to the Greek *daktylos* with the same meaning. The Greek word is also the origin of the name of a fruit, *date*. If you look at dates, not in a box, but growing high up in a palm *Phoenix dactylifera*, you will appreciate the similarity—*dactylifera* means ‘carrying fingers.’ So digital geometry is about counting on fingers and toes (perhaps implying using a system of base twenty) but if that seems to you not to be sweet enough, think of counting dates.

means that it is determined by its images on points, i.e., $F(A) = \bigcup_{x \in A} F(\{x\})$. So it is enough to know the digitization of an arbitrary point in X ; however, nothing requires the $F(\{x\})$ to be singleton sets.

The following setup seems to be sufficiently flexible.

Definition 10.1. *Let two sets X and Z be given, Z being a subset of X . Let there be given, for every $p \in Z$, a subset $C(p)$ of X , called the cell with nucleus p . Then the digitization determined by these cells is defined by*

$$(10.1) \quad F(\{x\}) = \{p \in Z; x \in C(p)\}$$

and

$$(10.2) \quad F(A) = \bigcup_{x \in A} F(\{x\}) = \{p \in Z; A \text{ meets } C(p)\}.$$

We may think of the cell $C(p)$ as a pixel or voxel, and of p as its address. If we think of C as a mapping $C: Z \rightarrow \mathcal{P}(X)$, then $F(A) = C^*(\mathcal{B}_A)$, where \mathcal{B}_A is the family of all cells which meet A .

It is clear that a digitization in this sense is always a dilation. As already pointed out, it is desirable that a nonempty set have a nonempty digitization; this is true if and only if the union of all cells is equal to the whole space X .

If X is an abelian group and Z a subgroup it is desirable that the digitization commute with translations, which means that $C(p) = C(0) + p$ for all $p \in Z$. Indeed, if $C(p) - p$ varies too much, it is easy to construct strange examples.

Example. A simple choice is $C(p) = \{p\}$. This yields the digitization $F(A) = Z \cap A$ already mentioned. If the set is fat, this digitization may work out well. In an abelian group with a metric we can even fatten the set using a dilation, defining $C(p)$ to be a ball $B_{\leq}(p, r)$ of radius r ; this yields $F(A) = Z \cap (A + B_{\leq}(0, r))$.

Example. If $X = \mathbf{R}$ and $Z = \mathbf{Z}$ we may choose $C(p) = [p - \frac{1}{2}, p + \frac{1}{2}]$. Then every set has a nonempty digitization, but the half-integers have a digitization consisting of two points. If we choose instead $C(p) =]p - \frac{1}{2}, p + \frac{1}{2}[$, then the digitization of a half-integer is empty. As a compromise we may choose $C(p) =]p - \frac{1}{2}, p + \frac{1}{2}]$; the digitization of a point is then always a point: $F(\{x\}) = \lceil x - \frac{1}{2} \rceil$. But then a new disadvantage appears: this digitization does not commute with the reflection $x \mapsto -x$.

Example. If $X = \mathbf{R}^2$ and $Z = \mathbf{Z}^2$ we may construct digitizations from what we have already done on the real axis. We may take

$$C(p) = [p_1 - \frac{1}{2}, p_1 + \frac{1}{2}] \times [p_2 - \frac{1}{2}, p_2 + \frac{1}{2}], \quad p \in \mathbf{Z}^2,$$

and similarly for the open and half-open intervals. Another choice is not to take the Cartesian product but to define $C(0)$, the cell with nucleus 0, as

$$(10.3) \quad C(0) = \{x; x_1 = 0 \text{ and } -\frac{1}{2} < x_2 \leq \frac{1}{2}\} \cup \{x; -\frac{1}{2} < x_1 \leq \frac{1}{2} \text{ and } x_2 = 0\},$$

and then $C(p)$ as $C(0) + p$. Thus $C(p)$ is a cross with center at p . This is the digitization used by Rosenfeld (1974). It is based on the mapping $\mathbf{R} \ni x \mapsto \lceil x - \frac{1}{2} \rceil \in \mathbf{Z}$, already

mentioned, a digitization of \mathbf{R} which takes a non-half-integer to the closest integer and moves down by one half in the case of half-integers. Let us call it the *Rosenfeld digitization* of \mathbf{R}^2 .

It is clear that in this case the union of the cells is very small compared with \mathbf{R}^2 , so that many sets have empty digitization. However, the union of all cells is equal to all grid lines $(\mathbf{R} \times \mathbf{Z}) \cup (\mathbf{Z} \times \mathbf{R})$, so that every straight line has a nonempty digitization. The same is true of a sufficiently long straight line segment. Thus this digitization can be used in the study of digital straight lines. Note that the family of all cells is disjoint, which implies that the digitization of a point is either empty or a singleton set.

The definition as such says nothing about how close a digitization of a point is to the point. To achieve this we must of course add some requirement that points in the cell $C(p)$ shall be reasonably close to p . This leads us to the next topic, that of Voronoi cells.

10.2. Voronoi cells

Let a metric space X be given as well as a subset Z . The metric of X shall be denoted by d . For a point $x \in X$ we view the points in Z close to x as approximations; there might be a best approximation. Given $p \in Z$ we shall give a name to the set of all x for which this particular p is a (not necessarily unique) best approximation: the *Voronoi cell*²⁶ with nucleus a is

$$\mathbf{Vo}(a) = \{x \in X; \forall b \in Z, d(x, a) \leq d(x, b)\}.$$

We also define the *strict Voronoi cell* as

$$\mathbf{Vo}_s(a) = \{x \in X; \forall b \in Z \setminus \{a\}, d(x, a) < d(x, b)\}.$$

Finally, one might define the *very strict Voronoi cell* as

$$\mathbf{Vo}_{vs}(a) = \{x \in X; d(x, a) < \inf_{b \in Z \setminus \{a\}} d(x, b)\}.$$

It is easy to construct examples where the very strict Voronoi cell is different from the strict Voronoi cell, but in all applications we are interested in they are equal.

Two different strict Voronoi cells are disjoint. Even more can be said: a (nonstrict) Voronoi cell is disjoint from every strict Voronoi cell with a different nucleus. The union of all strict Voronoi cells is almost equal to the whole space X ; there is only some garbage left out: these are the points which have at least two best approximants in Z . However, since we are mathematicians, we do not have the right to throw away that garbage; we must be careful and consider both the strict and the nonstrict Voronoi cells.

We now return to the topic of digitization. It seems reasonable that the digitization of a point should be contained in the set of all nuclei of Voronoi cells which contain that point. After all, these nuclei are the best approximants in Z of the point. This argument leads us to the following definition.

²⁶Named for Georgi Fedoseevič Voronoi (1868–1908)....

Definition 10.2. Let X be a metric space and Z a subset of X such that $Z \cap B_{<}(c, r)$ is finite for all $c \in X$ and all $r \in \mathbf{R}$. A Voronoi digitization of X into Z is a dilation $\text{Dig}: \mathcal{P}(X) \rightarrow \mathcal{P}(Z)$ such that

$$(10.4) \quad \text{Dig}(\{x\}) \subset \{a \in Z; x \in \text{Vo}(a)\}.$$

Note that if x belongs to some strict Voronoi cell $\text{Vo}_s(c)$, then it can belong to only one Voronoi cell, viz. the nonstrict cell $\text{Vo}(c)$ with the same nucleus, so that the right-hand side in (10.4) is a singleton set. Hence, $\text{Dig}(x)$ is either empty or equal to the singleton set $\{c\}$. But if x belongs to, say, two Voronoi cells, the right-hand side in (10.4) consists of a set $\{a, b\}$ with $a \neq b$, and there is a choice: $\text{Dig}(x)$ may be equal to \emptyset , $\{a\}$, $\{b\}$, or $\{a, b\}$. And if x belongs to m Voronoi cells, the value can be any of 2^m subsets of Z .

Thus $\text{Dig}(x)$ is either empty or a singleton set whenever x belongs to the union of all strict Voronoi cells, but in the complement, the value of the function may be a set with several elements. In some situations we do make a choice and define $\text{Dig}(x)$ to be a singleton set by introducing some new criterion. In fact, we have already done so when we defined the Khalimsky topology. If $X = \mathbf{R}$ and $Z = \mathbf{Z}$, then the Voronoi cells are the intervals $[n - \frac{1}{2}, n + \frac{1}{2}]$ and the strict cells are the open intervals $]n - \frac{1}{2}, n + \frac{1}{2}[$, $n \in \mathbf{Z}$. It is clear the digitization of a real number which is not of the form $n + \frac{1}{2}$ is the empty set or $\{\lfloor x + \frac{1}{2} \rfloor\}$. When $x = n + \frac{1}{2}$, we may choose $F(\{x\})$ to be \emptyset , $\{n\}$, $\{n + 1\}$, or $\{n, n + 1\}$. When we defined the Khalimsky topology, we chose $\{n\}$ for n even and $\{n + 1\}$ for n odd. But this is of course only one of many admissible choices.

Example. We get examples of Voronoi digitizations by taking $C(p) = \text{Vo}(p)$ or $C(p) = \text{Vo}_s(p)$. Sometimes it is possible to choose a cell in between these two so that the space is covered exactly once by the different cells; an example was already mentioned: if $X = \mathbf{R}^n$, $Z = \mathbf{Z}^n$ we may choose $C(p) = \prod]p_j - \frac{1}{2}, p_j + \frac{1}{2}]$.

Example. The digitization used by Rosenfeld (1974) is a Voronoi digitization, since the cell $C(p)$ defined in (10.3) is contained in the Voronoi cell, which is $\text{Vo}(p) = \{x \in \mathbf{R}^2; \|x - p\|_\infty \leq \frac{1}{2}\}$.

10.3. Digital lines

In \mathbf{R}^2 we know what a straight line is: it is a set of the form $\{(1 - t)a + tb; t \in \mathbf{R}\}$, where a and b are two distinct points in the plane. And a straight line segment is a connected subset of that line. We shall consider closed segments of finite length only, and may then write them as $\{(1 - t)a + tb; 0 \leq t \leq 1\}$, where a and b are the endpoints. We shall denote this segment by $[a, b]$.

We shall choose $Z = \mathbf{Z}^2$ in the discussion that follows. The digitization of a straight line segment is the image under Dig of $[a, b]$, thus

$$\text{Dig}([a, b]) = \bigcup_{t \in [0, 1]} F(\{(1 - t)a + tb\}) \subset \mathbf{Z}^2.$$

Suppose we are dealing with a Voronoi digitization. When $x = (1 - t)a + tb$ belongs to a strict Voronoi cell, which in this case is $\text{Vo}_s(p) = \{x; \|x - p\|_\infty < \frac{1}{2}\}$, $p \in \mathbf{Z}^2$, the function has the value $\{(\lfloor x_1 \frac{1}{2} \rfloor, \lfloor x_2 \frac{1}{2} \rfloor)\}$. When x_1 is a half-integer, and x_2 is not, the

digitization may be empty or consist of one or two points; when both coordinates are half-integers, the value may be a set of zero, one, two, three or four points.

In his famous paper (1974), Azriel Rosenfeld defined the digitization as in (10.3). In particular a point is always mapped to a point. For straight lines with slope less than 45° , he considered the intersections of its line segments with the vertical grid lines only. However, a line segment may intersect a horizontal grid line but no vertical grid line at all. In this case the cell is just the first segment in the union (10.3), but it does not matter so much, since the result will be trivially true for empty digitizations and the digitization is nonempty anyway for sufficiently long line segments.

We shall say with Rosenfeld that a subset A of \mathbf{R}^2 has *the chord property* if for all points $a, b \in A$ the segment $[a, b]$ is contained in $A + B_{<}(0, 1)$, the dilation of A by the open unit ball (or disk or square) for the l^∞ metric.

Theorems 10.3 and 10.4 are due to Rosenfeld (1974) and give together a characterization of the digitization of a straight line segment. The proof of Theorem 10.4 is new and is much shorter than the original proof.²⁷

Theorem 10.3. *The Rosenfeld digitization of a straight line segment has the chord property.*

Example. Let A be the set consisting of the five points $(0, 0), (1, 0), (2, 0), (3, 1), (4, 2)$. This set does not have the chord property. Indeed, the point $(2, 1)$ belongs to the segment $[(0, 0), (4, 2)]$, but it does not belong to the dilated set $A + B_{<}(0, 1)$, although it does belong to the closed set $A + B_{\leq}(0, 1)$. Thus, in view of the theorem, it cannot be the Rosenfeld digitization of a straight line segment. However, we may define a Voronoi digitization by declaring the digitization of $(0, -\frac{1}{2})$ to be $(0, 0)$, that of $(2, \frac{1}{2})$ to be $(2, 0)$, and that of $(4, 1\frac{1}{2})$ to be $(4, 2)$. Then A is the digitization of the straight line segment $[(0, -\frac{1}{2}), (4, 1\frac{1}{2})]$. This digitization does not commute with translations, which offers a kind of explanation—of course it should not be allowed to move up by one half from $(0, -\frac{1}{2})$ and $(4, \frac{3}{2})$ and down by one half from $(2, \frac{1}{2})$. Rosenfeld avoided this by always moving down in the case of half-integers.

Example. Slightly more generally we consider a set A consisting of five or six points $(0, 0), (1, 0), (2, 0)$, and $a = (a_1, a_2), (a_1 - 1, a_2 - 1), (a_1 - 2, a_2 - 2)$, where $a_1 \geq 4, a_2 \geq 2$. (If $a = (4, 2)$ we get the former example.) Then for no choice of a does this set have the chord property. Indeed, if $a_2 \geq \frac{1}{2}a_1$, then the point $(2, 2a_2/a_1)$, which is on the segment $[(0, 0), (a_1, a_2)]$, does not belong to $A + B_{<}(0, 1)$; if on the other hand $a_2 \leq \frac{1}{2}a_1$, then $(a_1 - 2, (a_1 - 2)a_2/a_1)$ on the same segment does not belong to $A + B_{<}(0, 1)$.

Proof of Theorem 10.3. Let L be a straight line and $D \subset \mathbf{Z}^2$ the digitization of a segment $[a, b]$ of L . Let p, q be two points in D and r an arbitrary point on the segment $[p, q]$. We shall prove that there exists a point $d \in D$ such that $\|d - r\|_\infty < 1$.

First we reduce to the case when the slope of L is between -1 and 1 —note that the whole setup is invariant under permutation of the coordinates. When it is exactly 0 or ± 1 the result is easy. When it is strictly between 0 and 1 we consider four cases: (1) $r_1 \in \mathbf{Z}$;

²⁷I found it during a refreshing walk in Hågadalen 2002-05-26.

- (2) $r_1 \notin \mathbf{Z}$ with $m < r_1 < m + 1$ for some $m \in \mathbf{Z}$;
 (2.1) L is under $[p, q]$ (in a non-strict sense) at $m + 1$;
 (2.2) L is over $[p, q]$ (in a non strict sense) at m ;
 (2.3) L is over $[p, q]$ at $m + 1$ and under $[p, q]$ at m .

In each of these cases we can find a point $d \in D$ with $\|d - r\|_\infty < 1$. When the slope is strictly between -1 and 0 things are similar.

In case (1) we take a point $s \in L$ with $s_1 = r_1$ and $|s_2 - r_2| \leq \frac{1}{2}$. Then there is a point $d \in D$ such that $d_1 = r_1 = s_1$ and $|d_2 - s_2| \leq \frac{1}{2}$, so $\|d - r\|_\infty \leq \|d - s\|_\infty + \|s - r\|_\infty \leq 1$. But can equality occur here? No. If we analyze the definition of the digitization we find that

$$d_2 - 1 < s_2 - \frac{1}{2} \leq r_2 < s_2 + \frac{1}{2} \leq d_2 + 1,$$

so that actually $|d_2 - r_2| < 1$, while $|d_1 - r_1| = 0$.

In case (2.1) we have $m < r_1 < m + 1$ and we know that L lies under $[p, q]$ at the abscissa $x_1 = m + 1$. We choose $s_1 = d_1 = m + 1$ and then find a point $d \in D$ with $|s_2 - r_2| \leq \frac{1}{2}$, $|s_2 - d_2| \leq \frac{1}{2}$. More precisely, denoting by r' the point on $[p, q]$ with abscissa equal to $m + 1$, we have $d_2 \geq s_2 - \frac{1}{2} \geq r'_2 - 1 > r_2 - 1$, for $r'_2 > r_2$; and $d_2 \leq s_2 + \frac{1}{2} \leq r'_2 + \frac{1}{2} < r_2 + 1$, for $s_2 \leq r'_2$ and the slope of L is strictly less than one half. This shows that $|d_2 - r_2| < 1$; we also know that $|d_1 - r_1| < 1$. Case (2.2) is similar to (2.1).

In case (2.3) we let s' be the intersection of L with the vertical line $x_1 = m$ and let s'' be the intersection of L with the line $x_1 = m + 1$. Thus $s', s'' \in L$ and $s'_1 = m$, $s''_1 = m + 1$. If L passes over $[p, q]$ at r we take $d \in D$ with $d_1 = m = s'_1$ and $|s'_2 - d_2| \leq \frac{1}{2}$. Then $d_2 > s'_2 - \frac{1}{2} > r_2 - 1$ and $d_2 \leq s'_2 + \frac{1}{2} < r_2 + \frac{1}{2}$, so that $|d_2 - s'_2| < 1$ and $d_2 = s'_2$. If on the other hand L is under $[p, q]$ at r , then we take $d \in D$ with $d_1 = m + 1 = s''_1$ and $|s''_2 - d_2| \leq \frac{1}{2}$. Then $d_2 \leq s''_2 + \frac{1}{2} < r_2 + 1$ and $d_2 > s''_2 - \frac{1}{2} \geq r_2 - 1$ so that $|d_2 - s''_2| < 1$ and $d_1 = s''_1$. The theorem is now completely proved.

To prove a converse we shall need the concept of *digital arc*, or *Khalimsky arc*. This is a continuous mapping from a finite Khalimsky interval into the Khalimsky plane \mathbf{Z}^2 . By abuse of language we shall call also the image of such a mapping a digital arc. It is thus a finite subset $\{p^1, \dots, p^m\}$ of \mathbf{Z}^2 such that p^j has exactly two adjacent points p^{j-1} and p^{j+1} for $j = 2, \dots, m - 1$, while p^1 has only one adjacent point p^2 , and p^m only one adjacent point p^{m-1} ; two points x and y are called *adjacent* if $\|x - y\|_\infty = 1$.

Theorem 10.4. *If a digital arc D in \mathbf{Z}^2 has the chord property, then it is the Rosenfeld digitization of some straight line segment in \mathbf{R}^2 .*

Lemma 10.5. *Denote by $\pi_j: \mathbf{Z}^2 \rightarrow \mathbf{Z}$ the projection $(x_1, x_2) \mapsto x_j$, $j = 1, 2$. If a digital arc D has the chord property, then one of the restrictions $\pi_j|_D: D \rightarrow \mathbf{Z}$, $j = 1, 2$, is injective.*

Proof. Assume that D is a digital arc and that, contrary to hypothesis, there are two pairs of points a, a' and b, b' in D such that $a_1 = a'_1$, $a_2 < a'_2$ and $b_2 = b'_2$, $b_1 < b'_1$. Then it is not difficult to see that D cannot have the chord property....

Proof of Theorem 10.4. Let D be a digital arc with the chord property. In view of the lemma and the symmetry of the digitization procedure, we may assume that there are no pairs of points a, a' in D with $a_1 = a'_1$, $a_2 \neq a'_2$. Given three real numbers α, β, γ

we define a strip in the plane by

$$S(\alpha, \beta, \gamma) = \{x \in \mathbf{R}^2; \alpha x_1 + \beta \leq x_2 \leq \alpha x_1 + \gamma\}.$$

Let us define the height of the strip as $\gamma - \beta$. The boundary $\partial S(\alpha, \beta, \gamma)$ of the strip has two components, given by the straight lines $x_2 = \alpha x_1 + \beta$ and $x_2 = \alpha x_1 + \gamma$. A set D of integer points is a subset of the digitization of a straight line segment if and only if D is contained in a strip of height strictly less than 1.

For every given α there is a smallest strip $S(\alpha, \beta, \gamma)$ containing D . Moreover, varying also α , there is a strip $S_0 = S(\alpha_0, \beta_0, \gamma_0)$ of smallest height. If D consists of only one or two points, the conclusion follows easily, so let us assume that D has at least three points. Clearly there must be at least one point of D in each component of the boundary of S_0 ; otherwise we could increase β or decrease γ to obtain a narrower strip. And one of these lines must contain a second point of D ; otherwise we could rotate the line slightly to obtain a strip of smaller height. For definiteness we shall assume that the three points on the boundary of the strip are p, s, q with $p_1 < s_1 < q_1$ and where p and q are on the lower boundary and s on the upper boundary. Let r be the point on $[p, q]$ with abscissa equal to that of s . (We note that p, s, q belong to \mathbf{Z}^2 , while r need not do so.) If D is not a subset of the digitization of a straight line, then the height of this smallest strip is at least 1, so that $s_2 \geq r_2 + 1$, showing that r does not belong to $s + B_{<}(0, 1)$. To see that D does not satisfy the chord property we must however show that there is no $d \in D$ such that $r \in d + B_{<}(0, 1)$. So far we only know that r does not belong to $s + B_{<}(0, 1)$. However, s is the only point in D on the vertical line $x_1 = s_1$ and all other points $d \in D$ satisfy $|d_1 - r_1| = |d_1 - s_1| \geq 1$, so that $\|r - d\|_\infty \geq |r_1 - d_1| \geq 1$. Therefore D does not satisfy the chord property.

We have thus proved that a digital arc D having the chord property is a subset of the digitization of some straight line L . However, since D is a digital arc, D must be connected. It is therefore the digitization of a connected subset of L . Obviously this subset can be taken to be compact, i.e., a straight line segment.

Generalize to other Voronoi digitizations....

Generalize to digitizations of convex sets....

10.4. Exercises

10.1. Let $X = \mathbf{R}^2$, $Z = \{(0, 0), (a_1, a_2)\}$ and determine the Voronoi cells

(a) when the metric is the Euclidean metric l^2 determined by the norm $\|\cdot\|_2$;

(b) when the metric is the l^∞ metric;

(c) when the metric is the l^1 metric.

10.2. Determine the Voronoi cells when $Z \subset \mathbf{C}$ is the set of all complex numbers $m + n\omega$, $m, n \in \mathbf{Z}$, where $\omega = \frac{1}{2} + \frac{i}{2}\sqrt{3}$ and we use the l^2 metric. What about other metrics?

References

Alexandroff, Paul (Aleksandrov, P. S.)

1935 Sur les espaces discrets. *C. R. Acad. Sci. Paris* **200**, 1649–1651.

1937 Diskrete Räume. *Mat. Sb.* **2** (44), 501–519.

Baclawski, Kenneth; Björner, Anders

- 1979 Fixed points in partially ordered sets. *Advances in Mathematics* **31**, 263–287.
- Birkhoff, Garrett
1940 *Lattice Theory*. New York City: American Mathematical Society. (Revised Edition 1948.)
- Borgefors, Gunilla
1984 Distance transformations in arbitrary dimensions. *Comput. Vision Graphics Image Process.* **27**, 321–345.
1986 Distance transformations in digital images. *Comput. Vision Graphics Image Process.* **34**, 344–371.
1994 Applications using distance transforms. In: *Aspects of Visual Form Processing*, pp. 83–108, (C. Arcelli; S. di Baja, Eds.), World Scientific, Singapore.
1996 On digital distance transforms in three dimensions. *Computer Vision and Image Understanding* **64**, 368–376.
- Bourbaki, Nicolas
1961 *Topologie générale*. Éléments de mathématique, première partie, livre III, chapitres 1 & 2. Third edition. Paris: Hermann.
1963 *Théorie des ensembles*. Éléments de mathématique, première partie, livre I, chapitre 3. Second edition. Paris: Hermann.
- Das, P. P.; Chatterji, B. M.
1988 Knight's distance in digital geometry. *Pattern Recognition Letters* **7**, 215–226.
- Everett, C. J.
1944 Closure operators and Galois theory in lattices. *Trans. Amer. Math. Soc.* **55**, 514–525.
- Gonzalez, Rafael C.; Woods, Richard E.
1993 *Digital Image Processing*. Addison-Wesley Publishing Company. xvi + 716 pp.
- Ghosh, Pijush K.; Kumar, K. Vinod
1998 Support function representation of convex bodies, its application in geometric computing, and some related representations. *Computer Vision and Image Understanding* **72**, 379–403.
- Halimskiĭ, E. D. (Efim Khalimsky)
1970 Applications of connected ordered topological spaces in topology. Conference of Math. Departments of Povolsia.
1977 *Uporyadochennnye topologicheskie prostranstva*. Kiev: Naukova Dumka. 92 pp.
- Herman, Gabor T.
1998 *Geometry of Digital Spaces*. Birkhäuser. x + 216 pp.
- Hilditch, J.; Rutovitz, D.
1969 Chromosome recognition. *Annals of the New York Academy of Sciences* **157**, 339–364.
- Hiriart-Urruty, Jean-Baptiste; Lemaréchal, Claude
1993 *Convex Analysis and Minimization Algorithms I. Fundamentals*. Springer-Verlag, vxii + 417 pp.
- Khalimsky, Efim; Kopperman, Ralph; Meyer Paul R.
1990 Computer graphics and connected topologies on finite ordered sets. *Topol. Appl.* **36**, 1–17.
- Kiselman, Christer O.
1969 Prolongement des solutions d'une équation aux dérivées partielles à coefficients constants. *Bull. Soc. Math. France* **97**, 329–356.

- 1996 Regularity properties of distance transformations in image analysis. *Computer Vision and Image Understanding*, **64**, No. 3, 390–398.
- 2000 Digital Jordan curve theorems. *Discrete Geometry for Computer Imagery*, 9th International Conference, DGCI 2000, Uppsala, Sweden, December 13–15, 2000. (Eds. Gunilla Borgefors, Ingela Nyström, Gabriella Sanniti di Baja.) Lecture Notes in Computer Science **1953**, pp. 46–56. Springer.
- MS Kan prickarna på datorskärmen bilda en kurva? Digital geometri ger svaret. [Can the dots on the computer screen form a curve? Digital geometry gives the answer. In Swedish.] Manuscript intended for the forthcoming book *Matematikens rikedomar*. 7 pp.
- Kong, T. Y.; Rosenfeld, A.
1989 Digital topology: Introduction and Survey. *Computer vision, graphics, and image processing* **48**, 357–393.
- Kong, Yung; Kopperman, Ralph; Meyer, Paul R.
1991 A topological approach to digital topology. *Amer. Math. Monthly* **98**, 901–917.
- Kopperman, R. D.; Kronheimer, E. H.; Wilson, R. G.
1998 Topologies on totally ordered sets. *Topol. Appl.* **90**, 165–185.
- Kronheimer, E. H.
1992 The topology of digital images. *Topol. Appl.* **46**, 279–303.
- Kuroš, A. G.
1962 *Lekcii po obščej algebre*. Moscow: F.M.
- Marchand-Maillet, Stéphane; Sharaiha, Yazid M.
2000 *Binary Digital Image Processing*. Academic Press. xxv + 251 pp.
- Matheron, G.
1967 *Éléments pour une théorie des milieux poreux*. Paris: Masson et C^{ie}.
1975 *Random sets and Integral Geometry*. New York: John Wiley & Sons. xxiii + 261 pp.
- Montanari, U.
1968 A method for obtaining skeletons using a quasi-Euclidean distance. *J. Assoc. Comput. Mach.* **15**, 600–624.
- Moore, Eliakim Hastings
1910 *Introduction to a Form of General Analysis*. New Have: Yale University Press.
- Moreau, J.-J.
1970 Inf-convolution, sous-additivité, convexité des fonctions numériques, *J. Math. Pures et Appl.* **49**, 109–154.
- Ore, Oystein
1944 Galois connexions. *Trans. Amer. Math. Soc.* **55**, 493–513.
- Pfaltz, J. L.; Rosenfeld, A.
1967 Computer representation of planar regions by their skeletons. *Comm. ACM* **10**, 119–125.
- Ragnemalm, I.,
1993 The Euclidean distance transform in arbitrary dimensions. *Pattern Recognition Letters* **14**, 883–888.
- Rockafellar, R. Tyrrell
1970 *Convex Analysis*. Princeton, NJ: Princeton University Press.
- Rosenfeld, Azriel

- 1974 Digital straight line segments. *IEEE Transactions on Computers*, **c-23**, No. 12, 1264–1269.
- 1979 Digital topology. *Amer. Math. Monthly* **86**, 621–630.
- Rosenfeld, A.; Pfaltz, J. L.
- 1966 Sequential operations in digital picture processing. *Journal of the ACM* **13**, 471–494.
- 1968 Distance functions in digital pictures. *Pattern Recognit.* **1**, 33–61.
- Serra, Jean
- 1982 *Image Analysis and Mathematical Morphology*. Academic Press. xviii + 610 pp.
- 1998 Connectivity on complete lattices. *Journal of Mathematical Imaging and Vision*. **9** (3), 231–251.
- 2001 *Lecture Notes on Morphological Operators*. Institut Mittag-Leffler, Lecture Notes No. 1, 2000/2001.
- Singer, Ivan
- 1997 *Abstract Convex Analysis*. New York: John Wiley and Sons, Inc. xxii + 491 pp.
- Starovoitov, Valery
- 1995 Toward a distance transform generalization. In: *Proceedings of the 9th Scandinavian Conference on Image Analysis*, pp. 499–506, (G. Borgefors, Ed.). Uppsala.
- Strömberg, Thomas
- 1996 The operation of infimal convolution. *Dissertationes Math.* **352**. 58 pp.
- Tarski, Alfred
- 1955 A lattice-theoretical fixpoint theorem and its applications. *Pacific J. Math.* **5**, 285–309.
- Thiel, Edouard
- 2001 *Géométrie des distances de chanfrein*. Mémoire scientifique; Habilitation à diriger des recherches. Marseille: Université de la Méditerranée (Aix-Marseille II). 100 pp.
- Tucker, A. W.
- 1936 Cell spaces. *Ann. of Math.* **37**, 92–100.
- Verwer, Ben J. H.
- 1991 Local distances for distance transformations in two and three dimensions. *Pattern Recognition Letters*, **12**, 671–682.
- Wyse, Frank, et al.
- 1970 Solution to problem 5712. *Amer. Math. Monthly* **77**, 1119.

Author's address: Uppsala University, Department of Mathematics,
P. O. Box 480, SE-751 06 Uppsala, Sweden.

E-mail: kiselman@math.uu.se

URL: <http://www.math.uu.se/~kiselman>

Telephone: +46-18-4713216 (office); +46-18-300708 (home)

Fax: +46-18-4713201