

7 Adjunctions

In which we use detailed examples and constructions from a variety of areas to illustrate the last of the really fundamental notions in category theory—that of adjunctions, or adjoint functors—and explore their key features.

In chapter 2, we saw how the notion of an isomorphism of categories—using a pair of functors $F: \mathbf{C} \rightarrow \mathbf{D}$ and $G: \mathbf{D} \rightarrow \mathbf{C}$ inverse to each other in the sense that $G \circ F = \text{id}_{\mathbf{C}}$ and $F \circ G = \text{id}_{\mathbf{D}}$ —was too restrictive to be useful, and that a far better category-theoretic notion of “sameness” of categories was supplied by the notion of an equivalence of categories, using natural transformations to relax the equalities and leave us with functors inverse “up to natural isomorphism,” in the sense that $G \circ F \cong \text{id}_{\mathbf{C}}$ and $F \circ G \cong \text{id}_{\mathbf{D}}$. The notion of an adjunction, or adjoint functors, is in a sense a further step in generalization, weakening the notion of an equivalence of categories, so that we are interested less in a relation (of sameness) between two categories and more in a relation between specific functors moving between those categories. Another perspective would be to say that adjunctions represent something like a further broadening of the notion of inverse, involving *unique* “reversal attempts” that supply the “closest thing” to inverses and capture important relations that exist even when inverses in the strict sense do not.

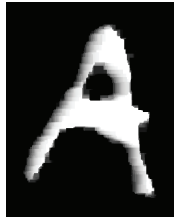
Like so many other important notions in category theory, these arise in an especially simple form in the special context of orders, so that adjoint relations between orders allow one to display the features of adjointness in a particularly accessible form. Adjoint functors between orders first appeared under the name of *Galois connections*. The next few examples will illustrate and motivate, via explicit examples in the context of orders, some of the many fundamental general features and properties of adjunctions.

7.1 Adjunctions through Morphology

Example 163 Suppose you receive a dark photocopy of some text, where the pen or marker appears to be bleeding:



With the help of your favorite programming language, you might perform what the image processing community would call an “erosion” of the image. After doing this (perhaps a few times), you would be left with something like



As one can see, erosion effectively acts to make thicker lines skinnier and detects, or enhances, the holes inside the letter A.

Suppose, instead, that you attempted a *dual* operation, called “dilation,” the effect of which is to thicken the image, so that lightly drawn figures are presented as if written with a thicker pen, and holes are (gradually) filled. In the case of dilating the original image, you would be left with something like



Now suppose that, for instance, after eroding the image you received, certain things have become harder to read. You decide that you would like to undo what you have done, perhaps because you have lost some important information. It seems sensible to hope you might undo it, and get back to the original image, by dilating the result of your erosion. But, in general, erosion and dilation do not admit inverses—in particular, they are not one another’s inverse—and there is no operation that would recover the exact original image from a dilated (or eroded) image. If an image is eroded and then dilated (or conversely), the resulting image will not be the original image. These operations *discard* information, so perhaps it is not so surprising that one would not get back to the original by “undoing” an erosion, for instance, by dilating the result.

However, in the failed search for an inverse to each operation, you will very quickly alight on something new, the basic properties of which appear to be useful and interesting in their own right. Erosion and dilation, while not inverses of each other, do seem to be related in a rather special way. In particular, eroding after we have dilated an image yields a very different result than dilating after eroding, even though neither composite gives back

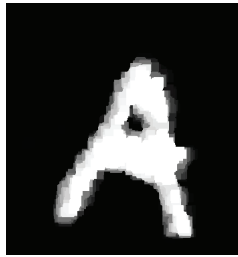
the original image. However, for an arbitrary image I , you will notice that it is “bounded” in both directions, in the sense of containing the result of one of the two composite operations while being contained by the result of the other, that is,

$$\text{Dilating after eroding } I \subseteq I \subseteq \text{Eroding after dilating } I.$$

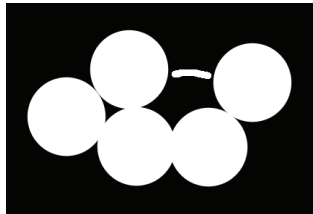
If we erode an image and then dilate the eroded image (making use of the same “structuring element” through both operations, on which more below), we arrive at a subset of the original image, in a process sometimes called *opening* the image by the image processing community. For instance, if we start with the following image



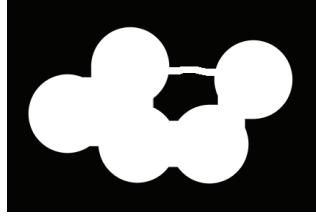
then opening it yields



Opening an image will leave an image that is generally smaller than the original, as it removes noise and protrusions and other small objects from an image, while preserving the shape and size of the more substantial objects in the image. On the other hand, dilating an image and then eroding the dilated image (with the same structuring element throughout)—sometimes called *closing* the image—leaves one with an image that is generally larger than the original. Starting with the following image



closing it will get rid of small holes, fill gaps in contours, smooth sections of contours, and fuse thin gulfs or breaks between figures. The result of closing the above yields



However, you may quickly learn that opening (or closing) an image twice leaves one with the same image as opening (or closing) it once. In other words, opening and closing are *idempotent* operations.

Altogether, the two basic operations of dilation and erosion are not quite inverses of one another, yet, as may already be evident from the discussion of the ways their idempotent composites “bound” the original above and below, they are nevertheless related in a special way, and are, in a sense, the “closest thing” to inverses (when these do not exist). We will explore that notion more closely and formally now.

Mathematical morphology is a field that deals with the processing of binary, gray-level, and other signals, and has proven useful in image processing. The majority of its tools are built on the two fundamental operators of dilation and erosion, and combinations thereof. It has a variety of applications involving image processing and feature extraction and recognition, including applications in x-ray angiography, biometrics, text restoration, among others.⁹⁹

A fundamental idea, in this setting, is that of “probing” an image with a basic, predefined shape and then examining how this fixed shape relates to the shapes comprising the image. One calls the probe the *structuring element*, which is itself a subset of the space, so that, for example, in the simple case of binary images, such a structuring element is itself just a binary image; these are, moreover, taken to have a defined origin. For instance, in the digital space $E = \mathbb{Z}^2$ (imagine a grid of squares), one might take for structuring element a 3×3 square, that is, the set

$$\{(-1, -1), (-1, 0), (-1, 1), (0, -1), (0, 0), (0, 1), (1, -1), (1, 0), (1, 1)\},$$

or perhaps a 3×1 rectangle with a designated central square, or a diamond or disk-shaped element, and so on. Going back to our earlier example of the blurry binary image letter A, we may take 0s to represent background and 1 for foreground, so that you basically have something like:

```

0 0 0 0 1 0 0 0 0
0 0 0 1 1 1 0 0 0
0 0 0 1 0 1 0 0 0
0 0 1 1 0 1 1 0 0
0 0 1 1 1 1 1 0 0
0 1 1 1 0 0 1 1 0
0 1 1 0 0 0 1 1 0
1 1 0 0 0 0 1 1 0
1 1 0 0 0 0 0 1 1
    
```

Then, we might take for structuring element $B \subseteq E$ a “diamond” (with origin boxed off)

```

    0  1  0
    1  1  1
    0  1  0
    
```

99. For classic introductions to mathematical morphology and filtering, see Serra (1986) and Serra and Vincent (1992).

In the simple case of our running example of a binary image in a bounded region, we took the structuring element B , placing B 's origin at each pixel of our image X as we scan over all of X , and compute at each pixel the dilation and erosion by taking the maximum or minimum value, respectively, of all pixels within the “window” or neighborhood covered by the structuring element (so that, e.g., in the case of dilation, a pixel is set to 1 if any of its neighboring pixels has the value 1).

More specifically, assuming we have fixed an origin in E , to each point p of E there will correspond the translation map that takes the origin to p ; such a map will then take B in particular onto B_p , the translate of B by p . In general, translation by p is a map $E \rightarrow E$ that takes x to $x + p$; thus, it takes any subset X of E to its *translate* by p ,

$$X_p = \{x + p \mid x \in X\}.$$

For a structuring element B , then, we can consider all its translates B_p . Given a subset (image) X of E , we can examine how the translates B_p of a given structuring element B interact with X . In the simple case of Boolean images (as subsets of a Euclidean or digital space), we carry out this examination via two operations:

$$X \oplus B = \{x + b \mid x \in X, b \in B\} = \bigcup_{x \in X} B_x = \bigcup_{b \in B} X_b,$$

called *Minkowski addition*, and its dual,

$$X \ominus B = \{p \in E \mid B_p \subseteq X\} = \bigcap_{b \in B} X_{-b}.$$

The former transformation (taking X into $X \oplus B$) is in fact what gives us a *dilation*, the basic property of which is that it distributes over union, while the latter (taking X into $X \ominus B$) is an *erosion*, the basic property of which is that it distributes over intersection. In the simple set-theoretical binary image case, dilations coincide with Minkowski addition; yet erosion of an image is the intersection of all translations by the points $-b$. In short,

Dilation of X by B is computed as the union of translations of X by the elements of B ,

while

Erosion of Y by B is the intersection of translations of Y by the reflected elements of B .

While we will see that we can give more general definitions of these operations, the particular behaviors of these operations in fact already follow from a general relationship underlying these operations.

Proposition 164 For every subset X, Y, B of our space E , where B is any structuring element, we have

$$X \oplus B \subseteq Y \text{ iff } X \subseteq Y \ominus B.$$

Proof. (\Rightarrow) Suppose $X \oplus B \subseteq Y$ and let $z \in X$ and $b \in B$. Then $z + b \in X \oplus B$, and thus $z + b \subseteq Y$. And $z + b \subseteq Y$ for any $b \in B$ implies that $z \in Y \ominus B$.

(\Leftarrow) Suppose $X \subseteq Y \ominus B$ and let $z \in X \oplus B$. Then there exists $x \in X$ and $b \in B$ such that $z = x + b$. But $x \in X$ and $X \subseteq Y \ominus B$ entails that $x \in Y \ominus B$. Thus, for every $b' \in Y$, we have $x + b' \in Y$, and in particular, $b \in B$, so $x + b \in Y$. But $z = x + b$, so $z \in Y$. \square

Before discussing the significance of this more generally, it is worth noting that there is no need to restrict attention, as we have thus far, to consideration of dilation and erosion in the simple case of Boolean images in digital space. We can also define the dilation and erosion of a *function* by a structuring element (itself regarded as a *structuring function*). For instance, given a function $f : E \rightarrow T$ from a space E to a set T of, for instance, gray-levels (i.e., a complete lattice that comes from a subset of $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$), and given a point $p \in E$, the *translate* of f by p is the function f_p whose graph is obtained by translating the graph $\{(x, f(x)) \mid x \in E\}$ by p in the first coordinate, that is, $\{(x + p, f(x)) \mid x \in E\}$, so that for all $y \in E$, we have

$$f_p(y) = f(y - p).$$

This defines the translation of a function by a point. Of course, if we extend this to the translation by a pair (p, t) , we get that for all $y \in E$,

$$f_{(p,t)}(y) = f(y - p) + t.$$

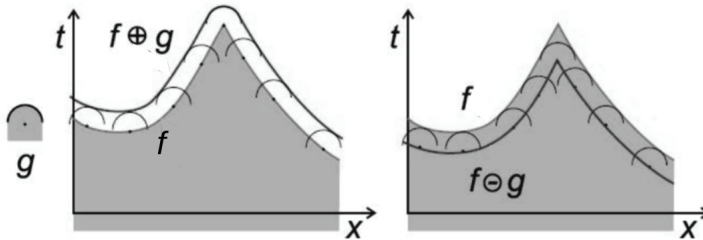
This approach allows us to define Minkowski addition and the dual operation for two *functions* $E \rightarrow T$. Where f plays the role of a (gray-level) image, and b is the functional analogue of a structuring element (i.e., a *structuring function*), for all $p \in E$, these operations will take on values

$$(f \oplus b)(p) = \sup_{y \in E} (f(y) + b(p - y))$$

and

$$(f \ominus b)(p) = \inf_{y \in E} (f(y) - b(y - p)).$$

Then the operator $\delta_g : T^E \rightarrow T^E$ taking $f \mapsto f \oplus g$ is *dilation* by g , and $\epsilon_g : T^E \rightarrow T^E$ taking $f \mapsto f \ominus g$ is *erosion* by g . In a low-dimensional case, using a portion of a disk for structuring function, these operators might do something like:¹⁰⁰



When we use a “flat structuring” element instead, such as that represented by a line, things are simplified even further, and we are effectively applying max and min filters, that is, for dilation

$$(f \oplus g)(x) = \sup_{y \in E, x-y \in B} f(y) = \sup_{y \in B_x} f(y)$$

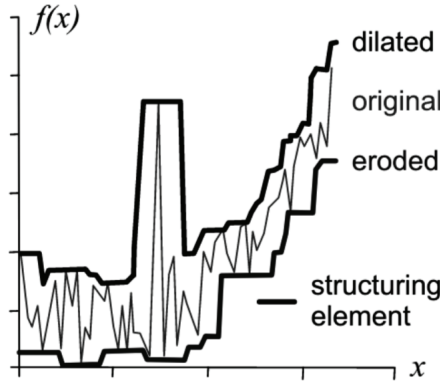
and for erosion

$$(f \ominus g)(x) = \inf_{y \in E, x-y \in B} f(y) = \inf_{y \in \check{B}_x} f(y),$$

where \check{B} is the transpose or symmetrical of B , that is, $\{-b \mid b \in B\}$. Erosion by a flat structuring function acts to shrink peaks and flatten valleys, while dilation acts in a dual fashion,

100. This image (and the one on the following page) is taken, with slight modifications, from Hlavac (2020).

flattening or rounding peaks and accentuating valleys. Taking, for instance, some price data, we might then have something like



Whether in its functional treatment, or in terms of the special relation $X \oplus B \subseteq Y$ iff $X \subseteq Y \ominus B$, the relation between dilation and erosion is part of a much more general and powerful story, exemplifying the notion of a Galois connection, itself an instance of the more general notion of an adjunction. We first supply the relevant definitions, and then further explore some of the powerful abstract features of this notion through the particular case of the operations of dilation and erosion.

Definition 165 Let $\mathcal{P} = (P, \leq_P)$ be a preordered set, and $\mathcal{Q} = (Q, \leq_Q)$ another preorder. Suppose we have a pair of monotone maps $F : P \rightarrow Q$ and $G : Q \rightarrow P$,

$$P \begin{array}{c} \xrightarrow{F} \\ \xleftarrow{G} \end{array} Q$$

such that for all $p \in P$ and $q \in Q$, we have the two-way rule

$$\frac{F(p) \leq_Q q}{p \leq_P G(q)}$$

where the bar indicates “iff.” If such a condition obtains, the pair (F, G) is said to form a monotone *Galois connection* between \mathcal{P} and \mathcal{Q} .

When such a connection obtains, we also say that F is the *left* (or *lower*) *adjoint* and G the *right* (or *upper*) *adjoint* of the pair, and write (for reasons we will see in a moment) $F \dashv G$ to indicate the relation.

Such a situation can be expressed in terms of the behavior of certain special arrows associated to each object of P and Q . In particular, let p be an object of P , and set $q = F(p)$. Then

$$\frac{F(p) \leq F(p)}{p \leq G(F(p))}$$

where the top is the identity arrow on $F(p)$ in Q , indicating reflexivity (which holds in any order). If we use θ to designate the bijection that realizes the “iff” situation, then $\theta(\text{id}_{F(p)})$ is a special arrow of Q , called the *unit* of p , where this arrow enjoys a certain universality property. There is a corresponding dual notion of a *counit*. In short, for each $p \in P$, we

call the *unit* an element $p \leq GF(p)$ that is least among all x with $p \leq G(x)$; dually, for each $q \in Q$, the *counit* is an element $FG(q) \leq q$ that is greatest among all y with $F(y) \leq q$. It can be shown that, given order-preserving maps $F : P \rightarrow Q$ and $G : Q \rightarrow P$, it is in fact *equivalent* to say (1) $F \dashv G$ and (2) $p \leq GF(p)$ and $FG(q) \leq q$.

Changing the variance of the functors involved gives us a slightly different notion.

Definition 166 Let $\mathcal{P} = (P, \leq_P)$ and $\mathcal{Q} = (Q, \leq_Q)$ be orders. Suppose we have a pair of antitone (order-reversing) maps $F : P \rightarrow Q$ and $G : Q \rightarrow P$,

$$P \begin{array}{c} \xrightarrow{F} \\ \xleftarrow{G} \end{array} Q$$

such that for all $p \in P$ and $q \in Q$, we have the two-way rule

$$\frac{q \leq_Q F(p)}{p \leq_P G(q)}$$

If such a condition obtains, the pair (F, G) is said to form an *antitone Galois connection* between \mathcal{P} and \mathcal{Q} .

We have seen many times now how any order \mathcal{P} can be regarded as a category by taking

$$x \leq_P y \text{ iff there exists an arrow } x \rightarrow y.$$

And in this setting, we further know that covariant functors between such categories are just monotone (order-preserving) functions, and that contravariant functors are antitone (order-reversing) functions. Thus, it is entirely natural to attempt to regard Galois connections in a more general, categorical guise. Doing so gives us the notion of an adjunction, which can accordingly be seen as a straightforward categorical generalization of the notion of a (monotone) Galois connection.

Definition 167 An *adjunction* is a pair of functors $F : \mathbf{C} \rightarrow \mathbf{D}$ and $G : \mathbf{D} \rightarrow \mathbf{C}$ such that there is an isomorphism

$$\text{Hom}_{\mathbf{D}}(F(c), d) \cong \text{Hom}_{\mathbf{C}}(c, G(d)),$$

for all $c \in \mathbf{C}, d \in \mathbf{D}$, which is moreover natural in both variables. When this obtains, we say F is left adjoint to G , or equivalently G is right adjoint to F , denoted $F \dashv G$.¹⁰¹

In saying that the isomorphism is “natural in both variables,” we mean that for any morphisms with domain and codomain as below, the square on the left commutes (in \mathbf{D}) iff the square on the right commutes (in \mathbf{C}):

$$\begin{array}{ccc} F(c) & \xrightarrow{f^\sharp} & d \\ F(h) \downarrow & & \downarrow k \\ F(c') & \xrightarrow{g^\sharp} & d' \end{array} \iff \begin{array}{ccc} c & \xrightarrow{f^\flat} & G(d) \\ h \downarrow & & \downarrow G(k) \\ c' & \xrightarrow{g^\flat} & G(d'). \end{array}$$

101. Sometimes the morphisms $F(c) \xrightarrow{f^\sharp} d$ and $c \xrightarrow{f^\flat} G(d)$ of the bijection given above are said to be *adjunct* or *transposes* of each other.

As one might expect, by considering functors of different variance, corresponding to *antitone* Galois connections, there is another notion, namely that of *mutually right adjoints* (and further, mutually left adjoints).

Definition 168 Given a pair of functors $F : \mathbf{C}^{op} \rightarrow \mathbf{D}$ and $G : \mathbf{D}^{op} \rightarrow \mathbf{C}$, if there exists a natural isomorphism

$$\text{Hom}_{\mathbf{D}}(F(c), d) \cong \text{Hom}_{\mathbf{C}}(G(d), c),$$

then we say that F and G are *mutually left adjoint*. Given the same functors, if there exists a natural isomorphism

$$\text{Hom}_{\mathbf{D}}(d, F(c)) \cong \text{Hom}_{\mathbf{C}}(c, G(d)),$$

then we say that F and G are *mutually right adjoint*.

Observe that an antitone Galois connection just names a mutual right adjoint situation between preorders (posets). Still following the example of the Galois connection definition, we can also recover the notions of the unit and counit of an adjunction.

Definition 169 Given an adjunction $F \dashv G$, there is a natural transformation

$$\eta : \text{id}_{\mathbf{C}} \Rightarrow GF,$$

called the *unit* of the adjunction. Its component

$$\eta_c : c \rightarrow GF(c)$$

at c is the transpose of the identity morphism $\text{id}_{F(c)}$.

Dually, there is a natural transformation $\mu : FG \Rightarrow \text{id}_{\mathbf{D}}$, called the *counit* of the adjunction, with component

$$\mu_d : FG(d) \rightarrow d$$

at d defined as the transpose of the identity morphism $\text{id}_{G(d)}$.

Any adjunction comes with a unit and a counit. In fact, conversely, given opposing functors $F : \mathbf{C} \rightarrow \mathbf{D}$ and $G : \mathbf{D} \rightarrow \mathbf{C}$, supposing they are equipped with natural transformations $\eta : \text{id}_{\mathbf{C}} \Rightarrow GF$ and $\mu : FG \Rightarrow \text{id}_{\mathbf{D}}$ satisfying a pair of conditions, then this data can be used to exhibit F and G as adjoint functors. In other words, we can use the natural transformations exemplifying the counit and unit maps, together with some conditions on these, to actually define an adjunction.

Definition 170 (*Adjunction, again*) An *adjunction* consists of a pair of functors $F : \mathbf{C} \rightarrow \mathbf{D}$ and $G : \mathbf{D} \rightarrow \mathbf{C}$, equipped with further natural transformations $\eta : \text{id}_{\mathbf{C}} \Rightarrow GF$ and $\mu : FG \Rightarrow \text{id}_{\mathbf{D}}$ satisfying what are sometimes called the *triangle identities*:

$$\begin{array}{ccc}
 F & \xrightarrow{F\eta} & FGF \\
 & \searrow \text{id}_F & \downarrow \mu F \\
 & & F
 \end{array}
 \qquad
 \begin{array}{ccc}
 G & \xrightarrow{\eta G} & GFG \\
 & \searrow \text{id}_G & \downarrow G\mu \\
 & & G.
 \end{array}$$

Then, the isomorphism $\text{Hom}_{\mathbf{D}}(F(c), d) \cong \text{Hom}_{\mathbf{C}}(c, G(d))$ realizing F and G as an adjoint pair, will exist precisely where there exists a pair of natural transformations, as above, satisfying the triangle identities.

Let us now return to the dilation and erosion of images and breathe some life into these ideas. The special relation that related Minkowski addition to its dual did not really depend on the particular form given to it in the translation-invariant case of a binary image, but exemplifies a more general notion of dilation and erosion on an arbitrary complete lattice, using the operations of supremum and infimum, where we have the adjunction

$$\text{dilate} \dashv \text{erode}.$$

To see how this works, we can first observe that both operations, dilation and erosion, are order-preserving (monotone), in the sense that $X \subseteq Y$ implies $X \oplus Z \subseteq Y \oplus Z$ and also $X \ominus Z \subseteq Y \ominus Z$. Moreover, while the order of an image intersection (union) and a dilation (erosion) cannot be interchanged freely, the dilation of the union of two images is indeed equal to the union of the dilations of the images, so the order can be interchanged; likewise, erosion of the intersection of two images yields the intersection of their erosions.

In dealing with the pair (dilate, erode), there is no need to be restricted to the poset of subsets of a digital space. There are clearly many choices for the underlying object space, where the images in question are held to reside. It is common to consider $\mathbb{P}(E)$ the space of all subsets of E (where E is the d -dimensional Euclidean space \mathbb{R}^d or the digital space \mathbb{Z}^d). But we might also consider $\text{Conv}(E)$ the space of all convex subsets of E ; or \mathcal{P}^E the space of “image functions” from E a discrete space to \mathcal{P} a pixel lattice; or the space $(T^3)^E$ of RGB color images (where RGB colors are triples (r, g, b) of numerical values, T^3 the lattice of RGB colors under componentwise order, and an RGB image is a function $E \rightarrow T^3$ taking each point $p \in E$ to a triple $(r(p), g(p), b(p))$ representing the RGB coloration of p); and so on. The takeaway, though, is that despite the differences in these underlying spaces, all such spaces form *complete lattices* (where this means that the underlying poset has all joins and meets). So if we consider a complete lattice \mathcal{L} , with the order \leq , supremum \bigvee , infimum \bigwedge , least element 0 , and greatest element I , such a lattice can be thought of as our “image lattice,” corresponding to a particular set of images we are working with. Traditionally, one then defines dilations and erosions as follows.

Definition 171 Let \mathcal{L} and \mathcal{M} denote complete lattices. For $\delta: \mathcal{L} \rightarrow \mathcal{M}$ and $\epsilon: \mathcal{M} \rightarrow \mathcal{L}$, we say that

- δ is a *dilation* provided for every $S \subseteq \mathcal{L}$,

$$\delta \left(\bigvee S \right) = \bigvee_{X \in S} \delta(X).$$

- ϵ is an *erosion* provided for every $T \subseteq \mathcal{M}$,

$$\epsilon \left(\bigwedge T \right) = \bigwedge_{Y \in T} \epsilon(Y).$$

Note that this also applies in the case of S, T empty, in which case a dilation is held to preserve 0 , while an erosion preserves I .

Our dilate-erode pair is actually an antitone Galois connection, where $\mathcal{M} = \mathcal{L}^{op}$, which just means that for all S, T in \mathcal{L}

$$\frac{T \leq^{op} \delta(S)}{S \leq \epsilon(T)},$$

which is, of course, the same as

$$\frac{T \geq \delta(S)}{S \leq \epsilon(T)}$$

or, equivalently,

$$\frac{\delta(S) \leq T}{S \leq \epsilon(T)},$$

where the order here is now the same, that given on \mathcal{L} , above and below the line. Thus, we have recovered the usual notion of an adjunction with $\delta: \mathcal{L} \rightarrow \mathcal{L}$ order-preserving and $\epsilon: \mathcal{L} \rightarrow \mathcal{L}$ order-preserving! Dilations and erosions are then precisely just the order-preserving (monotone) transformations on a complete lattice that moreover commute with the supremum and infimum, respectively.¹⁰²

Morphological operators are thereby given a unified treatment in the general framework of an adjoint pair on complete lattices. A number of well-established properties concerning the interaction of these operators then fall out immediately from the general framework of adjunctions. Conversely, we can illustrate such general facts via the present operators on images.

Suppose we have an adjunction $\delta \dashv \epsilon$ on a complete lattice \mathcal{L} .¹⁰³ Then a number of morphologically significant facts come “for free” as corollaries of general categorical truths about an adjoint pair. Even the fact that δ is a dilation and ϵ is an erosion in the first place can be derived from the existence of this special adjoint relationship. In what follows, we explore some of these general truths through the lens of some notable particular truths about dilations and erosions.

7.1.1 Uniqueness of Adjoints

Proposition 172 To each dilation δ there corresponds a unique erosion ϵ , namely

$$\epsilon(X) = \bigvee \{S \in \mathcal{L} \mid \delta(S) \leq X\},$$

and to each erosion ϵ there corresponds a unique dilation,

$$\delta(X) \bigwedge \{S \in \mathcal{L} \mid \epsilon(S) \geq X\}.$$

102. Note that, if we were to regard the pair as comprising an antitone Galois connection, then we would be saying that the operators exchanged suprema and infima, in the sense that, for example, $\delta(\bigvee_i x_i) = \bigwedge \delta(x_i)$.

103. Looking ahead to what will have to be true of such functors, they have been given the names δ and ϵ , suggesting dilations and erosions; however, at this point, we do not yet require or assume anything about the maps δ and ϵ , except that they form an adjoint pair moving between \mathcal{L} and itself.

This ultimately derives from a general result that assures us that, like inverses, adjoints are unique (well, actually “unique up to unique isomorphism,” but we can ignore this in our special case):

Proposition 173 Adjoint maps are unique.

In the case of orders, with order-preserving maps between them, this just means

1. if F_1 and F_2 are left adjoints of G , then $F_1 = F_2$;
2. if G_1 and G_2 are right adjoints of F , then $G_1 = G_2$.

Proof. (We focus on the simple case of orders, and prove (1); (2) follows by duality) From the adjointness assumptions, we have both

$$\frac{F_1(p) \leq q}{p \leq G(q)}$$

and

$$\frac{p \leq G(q)}{F_2(p) \leq q},$$

so immediately we have that $F_1(p) \leq q$ iff $F_2(p) \leq q$. Set $q = F_1(p)$, making $F_1(p) \leq q$ trivially true, forcing $F_2(p) \leq F_1(p)$ to be true as well. Similarly, set $q = F_2(p)$ and use the trivial truth $F_2(p) \leq F_2(p)$ to force $F_1(p) \leq F_2(p)$. In a poset, this entails that $F_1(p) = F_2(p)$, p arbitrary. \square

The adjunction then gives rise to the formulas

$$G(q) = \bigvee \{p \mid F(p) \leq q\}$$

and

$$F(p) = \bigwedge \{q \mid p \leq G(q)\},$$

which displays the uniqueness of the adjoints, and so explains the unique erosion (dilation) corresponding to each dilation (erosion), as written above.

In general, a given map may or may not have a left (or right) adjoint; the map may have one without the other, neither, or both (where these may be the same or different). But if it does have a left (or right) adjoint, we can be confident that, even though they are not quite inverses, the adjoint is unique up to isomorphism.

Adjoint functors also interact in particularly interesting and useful ways with the limit and colimit constructions, a connection we now explore.

7.1.2 Limit and Colimit Preservation

Proposition 174 δ is a dilation and ϵ is an erosion, and both are order-preserving.

This follows immediately from a very important category-theoretic result, namely that

Proposition 175 Right adjoints preserve limits (*RAPL*); left adjoints preserve colimits (*LAPC*).¹⁰⁴

104. Terminologically, recall that a general functor that is limit-preserving is said to be a *continuous* functor, while a colimit-preserving functor is a *cocontinuous* functor. Another related concept we will make use of later in the book is the following: a functor is said to be *left exact* if it preserves *finite* limits, and *right exact* if it preserves *finite* colimits. Speaking of (co)limits, it is worthwhile noting that entities exhibiting universality, like colimits and limits, can themselves be phrased entirely in terms of adjoint functors. Then, one of the advantages of this

Instead of proving this in the general case, we will show how it obtains in our special case of maps between orders, in which setting limits are infima (meets) and colimits are suprema (joins).

Proposition 176 *RAPL and LAPC in the special case of orders:*

1. If $f: \mathbb{Q} \rightarrow \mathcal{P}$ has a right adjoint (i.e., is a left adjoint), then it preserves the suprema that exist in \mathbb{Q} .
2. If $g: \mathcal{P} \rightarrow \mathbb{Q}$ has a left adjoint (i.e., is a right adjoint), then it preserves the infima that exist in \mathcal{P} .

Proof. (Of (1), since (2) follows by duality) Assume $S = \{q_i\}_{i \in I}$ is a family of elements of \mathbb{Q} with a supremum $\bigvee S$ in \mathbb{Q} . Claim: $f(\bigvee S)$ is the supremum in \mathcal{P} of the family $\{f(q_i)\}_{i \in I}$, that is,

$$f\left(\bigvee S\right) = \bigvee f(S).$$

But $f(\bigvee S) \leq p$ iff $\bigvee S \leq g(p)$ (since, by assumption, f has a right adjoint, call it g). And this latter inequality holds iff for all $q_i \in S$, we have $q_i \leq g(p)$. But then we can again use the assumed adjoint relation $f \dashv g$, and see that this latter inequality will hold iff $f(q_i) \leq p$ for all $q_i \in S$, and this in turn will hold iff for all $t \in f(S)$, we have $t \leq p$. In sum, then, we have that $f(\bigvee_i q_i) \leq p$ if and only if $\bigvee_i f(q_i) \leq p$, or that f preserves any suprema that exist in \mathbb{Q} . □

But a dilation (erosion) was just *defined* as an order-preserving map that commutes with colimits (limits). So δ being a left adjoint suffices to tell us that δ must be a dilation (the dual situation holding for an erosion).

7.1.3 Adjoints Compose

Proposition 177 Given two dilations $\delta: \mathcal{L} \rightarrow \mathcal{M}, \delta': \mathcal{M} \rightarrow \mathcal{N}$ and two erosions $\epsilon: \mathcal{M} \rightarrow \mathcal{L}, \epsilon': \mathcal{N} \rightarrow \mathcal{M}$ such that $\delta \dashv \epsilon$ and $\delta' \dashv \epsilon'$, then their composition forms an adjunction $\delta' \circ \delta \dashv \epsilon \circ \epsilon'$.

This exemplifies a general result in category theory, namely:

Proposition 178 Left (right) adjoints are closed under composition, that is, given the adjunctions

$$\mathbf{C} \begin{array}{c} \xrightarrow{F} \\ \dashv \\ \xleftarrow{G} \end{array} \mathbf{D} \begin{array}{c} \xrightarrow{F'} \\ \dashv \\ \xleftarrow{G'} \end{array} \mathbf{E},$$

the composite $F' \circ F$ is left adjoint to the composite $G \circ G'$:

$$\mathbf{C} \begin{array}{c} \xrightarrow{F' \circ F} \\ \dashv \\ \xleftarrow{G \circ G'} \end{array} \mathbf{E}.$$

In this way, arbitrarily long strings of adjoints can be produced.

adjunction perspective is that the (co)limit of every **J**-shaped diagram in **C** can be defined all at once, rather than just taking the (co)limit of a particular **J**-shaped diagram $X: \mathbf{J} \rightarrow \mathbf{C}$. We will exhibit the relevant adjunction in the final section of this chapter.

Moreover, another fact from morphology follows from the facts that adjoints compose (and using LAPC and RAPL), namely that for dilations and erosions on the same complete lattice, if $\delta_j \dashv \epsilon_j$ forms an adjoint pair for every $j \in J$, then $(\bigvee_j \delta_j, \bigwedge_j \epsilon_j)$ is an adjunction.

7.1.4 Units and Counits

The “opening” operator bounds an image on the left, while its “closing” bounds it on the right, that is,

Proposition 179

$$\delta\epsilon \leq \text{id} \leq \epsilon\delta.$$

This is immediate from the unit natural transformation, $\text{id} \leq \epsilon\delta$, and counit $\delta\epsilon \leq \text{id}$.

7.1.5 Fixed Point Formulae

Proposition 180 $\delta\epsilon\delta = \delta$ and $\epsilon\delta\epsilon = \epsilon$.

The unit and counit maps, satisfying the triangle identities, give the following general “fixed point formulae” result underlying the above:

Proposition 181 If \mathcal{P} and \mathcal{Q} are posets and $F : \mathcal{P} \rightarrow \mathcal{Q}$ and $G : \mathcal{Q} \rightarrow \mathcal{P}$ form a (monotone) Galois connection (adjunction), with $F \dashv G$, then the following *fixed point formulae* will hold for F and G :

$$FGF = F \text{ and } GFG = G.$$

Proof. The triangle identities give $F(p) \leq FGF(p) \leq F(p)$ for all $p \in \mathcal{P}$, so $F = FGF$. The second formula follows similarly. \square

7.1.6 Returning to Main Discussion

In the last two items, we saw how the unit and counit maps determine two important endomaps, namely $\delta \circ \epsilon$ (“opening”) and $\epsilon \circ \delta$ (“closing”). The presence of unit and counit further give us the fixed point formulae, which translates to the morphologically significant fact,

$$\delta\epsilon\delta = \delta \quad \text{and} \quad \epsilon\delta\epsilon = \epsilon.$$

This “stability” property of openings and closings means, in terms of the interpretation of such operations as filters, that they effectively “complete their task” (unlike many other filters, where repeated applications can involve further modifications of the image, with no guarantee of the outcome after a finite number of iterations). In general, the above fixed point formula further entails, in particular, that $\epsilon\delta$ and $\delta\epsilon$ are each idempotent. Thus, altogether, the composite monotone map $\epsilon\delta$, for its part, has the properties that

- $p \leq \epsilon\delta(p)$, and
- $\epsilon\delta\epsilon\delta(p) = \epsilon\delta(p)$.

But this is exactly to say that $\epsilon\delta$ is a *closure* operator, in the following general sense.

Definition 182 A *closure operator* on a poset \mathcal{P} (typically some poset of subobjects, for example, the powerset poset)¹⁰⁵ is an endomap $K : \mathcal{P} \rightarrow \mathcal{P}$ such that

1. for each $p \leq p' \in \mathcal{P}$, $K(p) \leq K(p')$ (monotonicity);
2. for each $p \in \mathcal{P}$, $p \leq K(p)$ (extensivity);
3. for each $p \in \mathcal{P}$, $K(K(p)) = K(p)$ (idempotence).

There is an important dual notion to closure, called the kernel operator (or dual closure), where this is an endomap that, like K , arises from a Galois connection, and is both monotone and idempotent, yet satisfies the dual of the extensivity property.

Definition 183 A *kernel operator* (or *dual closure*) is an endomap L satisfying

1. for each $p \leq p' \in \mathcal{P}$, $L(p) \leq L(p')$ (monotonicity);
2. for each $p \in \mathcal{P}$, $L(p) \leq p$ (contractivity);
3. for each $p \in \mathcal{P}$, $L(L(p)) = L(p)$ (idempotence).

In short, these notions are all part of a much more general story, namely that for a Galois connection or adjunction on posets such as $\delta \dashv \epsilon$ as above, the composite $\epsilon \circ \delta$ will automatically be monotone, extensive, and idempotent, that is, a closure operator on the underlying poset (or lattice) \mathcal{P} ; dually, $\delta \circ \epsilon$ will be monotone, contracting, and idempotent, that is, a kernel operator on \mathcal{Q} .

Before leaving this example, we will explore a few last notions via morphology. We will let the induced kernel operator—which is precisely what the morphology community calls by the name “opening”—be denoted $\phi = \delta\epsilon$, while $\kappa = \epsilon\delta$ will denote the induced closure (or “closing” operator, to be consistent with the mathematical morphology literature). In the binary case, opening and closing are typically defined, respectively, as

$$X \circ B = (X \ominus B) \oplus B = \bigcup \{B_p \mid p \in E, B_p \subseteq X\}$$

$$X \bullet B = (X \oplus B) \ominus B.$$

Morphological closing is just dilation (by some B) followed by erosion of the result by B , while morphological opening is the erosion (by some B) followed by dilation of the resulting image by B . Closing acts to fill out narrow holes. In terms of translations with the structuring element, the opening of an image A by B is the complement of the union of all translations of B that fall outside (do not overlap) A . As extensive (i.e., larger than the identity mapping), closings of an image are generally “larger” than the original image. Opening, for its part, acts to remove noise, narrow connections between regions, and parts of objects, generally attenuating peaks and other small protrusions or components. If you have a note where the writing appears to be growing tiny roots from its edges, opening effectively acts to remove these outer leaks at the boundary, rounding the edges. In terms of translations with the structuring element, the opening of A by B is the union of all translations of B that fit completely within A . As antiextensive (contracting), openings of an image are generally “smaller” than the original.

105. “Subobjects” are formally introduced in chapter 11. For now, you can just think of them as categorical generalizations of the notion of a *subset*.

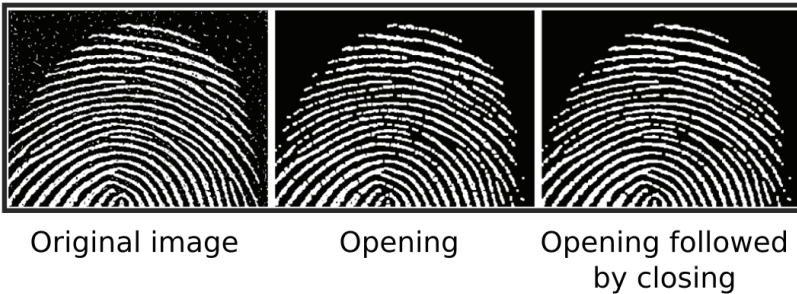
Exercise 19 Composing dilations and erosions, we found the composite operations of opening ($\phi = \delta\epsilon$) and closing ($\kappa = \epsilon\delta$), which were moreover idempotent. Further composing openings and closings with one another (e.g., $\kappa \circ \phi$), how many more distinct operations can we produce? Describe, in terms of their effect on images, at least one of these “image filters.” Finally, consider how the composite operators must be related to one another.

Solution By (alternately) composing openings $\phi (= \delta\epsilon)$ and closings $\kappa (= \epsilon\delta)$, we can obtain four new filters in total, each four of which are idempotent.

1. closing-after-opening: $\kappa\phi$
2. opening-after-closing: $\phi\kappa$
3. opening-after-closing-after-opening: $\phi\kappa\phi$
4. closing-after-opening-after-closing: $\kappa\phi\kappa$.

You can easily convince yourself that no other nontrivial operator can be obtained by further composition with any combination of ϕ s and κ s. Any attempt to produce further new operations by pre- or postcomposing the above four with κ or ϕ will just reduce back to one of those four, by the idempotence of these operators (together with the idempotence of κ and ϕ themselves).

These composites are used in the course of various image processing tasks, such as “smoothing” an image or performing image segmentation. An example of $\kappa\phi$ is given by the following:¹⁰⁶



In terms of the relations between these four (together with the original opening and closing operators as well), it is easy to show that

$$\phi \leq \phi\kappa\phi \leq \left\{ \begin{matrix} \kappa\phi \\ \phi\kappa \end{matrix} \right\} \leq \kappa\phi\kappa \leq \kappa,$$

and moreover $\phi\kappa\phi$ will be the greatest filter smaller than $\phi\kappa \wedge \kappa\phi$, while $\kappa\phi\kappa$ will be the smallest filter greater than $\phi\kappa \vee \kappa\phi$.

7.2 Adjunctions through Modalities

We can get an even better handle on adjunctions by looking at further applications of such notions, specifically some applications involving *modalities*. At least since Aristotle’s attempt to understand certain statements containing the words “necessary” and “possible,” philosophers and logicians have been interested in the logic of different operators

106. This image is taken from Bobick (2014).

describing different *ways of being true*. Modal logic began as the study of *necessary* and *possible* truths, but over the course of at least the last 100 years it has been recognized that modalities abound in both natural and formal languages. These days modal logic is more commonly regarded as the much broader study of a variety of constructions that modify the truth conditions of statements (which includes, most notably, statements concerning knowledge, belief, ethics, temporal happenings, how computer programs behave).¹⁰⁷ Moreover, one can construe modal operators in terms of adjunctions—and, in this way, modalities can be shown to arise in a number of other more unexpected settings. In particular, there are in fact a number of close connections between erosion and dilation, and the modal operators \Box of necessity and \Diamond of possibility, respectively.

After a brief subsection devoted to establishing some background on the more traditional logical treatment of modalities, the subsequent four subsections realize all these ideas in four different settings: (1) in connection with the treatment of *negation*, (2) in the **Qua** category, (3) in graphs, and (4) in topology.

7.2.1 Some Background on Modalities

The reader is likely already familiar with classical propositional logic (PL). Classical PL is ultimately built from propositional variables p, q, r, \dots , where these are variables (representing *propositions*) that can be assigned a truth value, and certain symbols called logical connectives, where this includes \neg (“not”), \wedge (“and”), \vee (“or”), \rightarrow (“implies”), \leftrightarrow (“if and only if”). The connectives enable compound propositions or formulas to be constructed from simpler propositional variables, and in such a way that the truth-value of the compound formula is defined as a function of the truth values of the simpler propositions. We define a *formula* by specifying that it is any expression constructed according to the following rules:

1. every propositional variable is a formula;
2. given any formulas ϕ and ψ already constructed, the expressions $\neg\phi$ (negation), the conjunction $\phi \wedge \psi$, the disjunction $\phi \vee \psi$, the implication $\phi \rightarrow \psi$ (“ ϕ implies ψ ”), and the biconditional $\phi \leftrightarrow \psi$ (“ ϕ if and only if ψ ”) are formulas.

Because one can combine simple propositions into compound ones in many ways, additional symbols like parentheses are needed to avoid ambiguities. In classical PL, an *interpretation* of a formula is a functional assignment of truth-values to its constituent variables. It is of course possible that a formula will have different truth-values under different interpretations—but a formula is said to be *valid* if it is true under every interpretation.

Intuitionistic PL is defined by the following axiom schemata (meaning one can substitute arbitrary formulae, obtaining *instances* of the axioms):

1. $\phi \rightarrow (\psi \rightarrow \phi)$
2. $(\phi \rightarrow (\psi \rightarrow \rho)) \rightarrow ((\phi \rightarrow \psi) \rightarrow (\phi \rightarrow \rho))$
3. $(\phi \wedge \psi) \rightarrow \phi$
4. $(\phi \wedge \psi) \rightarrow \psi$
5. $\phi \rightarrow (\psi \rightarrow (\phi \wedge \psi))$

107. For a nice history of formal approaches to modal logic, see Goldblatt (2003).

6. $\phi \rightarrow (\phi \vee \psi)$
7. $\psi \rightarrow (\phi \vee \psi)$
8. $(\phi \rightarrow \rho) \rightarrow ((\psi \rightarrow \rho) \rightarrow (\phi \vee \psi) \rightarrow \rho)$
9. $(\phi \rightarrow \psi) \rightarrow ((\phi \rightarrow \neg\psi) \rightarrow \neg\phi)$

together with one inference rule (*modus ponens*)

$$\frac{\phi \quad \phi \rightarrow \psi}{\psi}$$

where this means that if ϕ and $\phi \rightarrow \psi$ are derivable from the axioms, then so too is ψ .

We then use the axioms of the system, together with the inference rule, to produce proofs or derivations—where a proof is just a finite sequence of formulas, usually displayed vertically, such that each line of the proof is either an axiom instance of or is inferable from earlier lines using the inference rule.

When we add to this list of axioms a tenth axiom,

$$\phi \vee \neg\phi,$$

the law of excluded middle (or $\neg\neg\phi \rightarrow \phi$), we get classical PL.

In general, an axiom system is said to be *sound* if every formula that is derivable from the axiom system is valid. An axiom system is said to be *complete* if every valid formula can be derived from this axiom system. The ten axioms we gave above supply one of the sound and complete axiom systems for classical PL.

One can then “upgrade” PL by augmenting it with quantifiers “for all” (\forall) and “there exists/there is at least one” (\exists), where the role of these operators is of course to tell us *of how many* a proposition is true. The resulting logic—predicate logic—supplies us with a way of proving the validity of a number of perfectly valid arguments that are simply invalid when symbolized in the somewhat limited notation of PL. For instance, the venerated syllogism

1. All humans are mortal
2. Socrates is a human
- \therefore Socrates is mortal

can only be symbolized in PL in a manner that makes it invalid—for instance, as

1. ϕ
2. ψ
- $\therefore \rho$.

Yet, of course, the syllogism ultimately represents a valid argument, and by augmenting PL with quantifiers, we can show this. PL, and its extension to predicate or quantifier logic (by adding the quantifiers \forall and \exists), is of course rather useful in the formalization and analysis of many arguments. However, in the same way that the syllogism above could not be shown valid without a way of reasoning with quantifications of statements, there are further valid arguments that cannot be shown to be valid using just propositional (or predicate) logic. Take, for instance, the following argument:

If a new course on sheaves is to be offered next year, then proposal submissions must be made to the department before September. If proposal submissions are to be made to the department before September, then a departmental meeting must be had. A month’s notice must be given if a

departmental meeting is to be had. But it is already August. Because it is not possible to give a month's notice, it follows that it is not possible to offer a new course on sheaves next year.

Surely we would like to be able to show that such an argument is valid. But in order to show that such an argument is indeed valid, we need a way of reasoning with what is expressed by these various ways of *qualifying* how a proposition is true—in particular, capturing the notions of *must* and *possible* used above. The idea of reasoning with such qualifications of “ways of being true”—where these are called *modalities*, the most conspicuous and extensively studied of which have been “*it is necessarily the case that...*” and “*it is possibly the case that...*”—has been studied by philosophers for millennia, since at least the time of Aristotle, and reasoning with such qualifications was extensively discussed and debated in the Middle Ages. Roughly, a modality is just a phrase or concept that is applied to a given statement (ϕ) to create a new statement, where this latter statement now makes an assertion of the mode of truth of ϕ —where or how ϕ is true, when ϕ is true, under which circumstances ϕ may be true.

As mentioned, modal logic began as the study of *necessary* and *possible* truths. But even confining our attention to logicians, these days modal logic is more commonly regarded as the much broader study of a variety of constructions that modify the truth conditions of statements. Here is a partial sample of some prominent modalities considered by modal logicians:

- epistemic logic: it is known (to agent X) that
- doxastic logic: it is believed that
- deontic logic: it is obligatory that
- dynamic logic: after the program/computation terminates, the program enables that
- metalogic/provability logic: it is provable that
- tense logic: at all future times it is true that

In the formal analysis of modalities a key feature to emerge was that many modal operators came in dual pairs, similar to how the universal quantifier \forall and the existential quantifier \exists are dual and so interdefinable using negation

$$\forall x(\dots x \dots) := \neg \exists x \neg (\dots x \dots)$$

or

$$\exists x(\dots x \dots) := \neg \forall x \neg (\dots x \dots).$$

In a similar fashion, necessity (symbolized using \Box) acts as a sort of universal counterpart to the notion of possibility (symbolized using \Diamond), which plays the role of its existential dual, in the sense that

$$\Box := \neg \Diamond \neg$$

$$\Diamond := \neg \Box \neg,$$

reflecting the idea that to hold that “it is necessarily the case that ...” is equivalent to maintaining that “it is not possible that it is not the case that ...”. As these operators are given different readings, and put to different use, key dualities are captured—for instance,

- “it is obligatory that ...” vs. “it is permissible that ...”;
- “at all future times it is true that ...” vs. “at some future time it is true that ...”;
- “it is provable that ...” vs. “it is consistent (not provable that it is not the case) that ...”

Realizing such matters as a case of an adjoint relationship allows us to connect the story of modalities and such dualities to an even broader class of structures, beyond the confines of logic. The next four sections offer a glimpse into some of those further connections.

7.2.2 On What Is Not

Since at least the time of one of the first Western philosophical texts, attributed to Parmenides (around 500 BCE), the nature of *negation* has been on people's minds. This includes a number of issues, such as the following:

- Would a complete description of *what is* need to include any description of *what is not*? In other words, what is the “ontological status” of the negated entities or negative states of affairs?¹⁰⁸
- When is the negation of a negation (negated entity) the identity (the original entity)?

One might further motivate such concerns as follows. One might try to argue that, philosophically, holes, shadows, fissures, boundaries—and other such “negative” or derivative entities—seem somehow less real or fundamental than (or at least not to be on the same footing as) the “positive” objects that produce or surround or support them. At the very least, this sort of observation seems to have some validity in that it does seem somehow more difficult to supply identity criteria for holes, for instance, compared to ordinary material objects (for holes appear to be *made of nothing*), or even to speak of what holes are (what are the *parts* of a hole?).

Today, we are most accustomed to thinking of negation as a linguistic or logical operator on a language, where the operation leads from an expression to the contradictory expression. Typically, the expression in question is a proposition or a part of a proposition. But we might attempt to regard negation, more broadly, as an operation that can also take place on larger wholes, on entire structures or theories. Moreover, one might argue that, however one approached it, the “right” understanding (and description) of negation would need to capture, above all, the relation or dependence between what (the structure) is being negated and the result of this operation of negation.

Such a perspective on negation is arguably exemplified in the facts that (1) negation is a contravariant functor on a particular category (to itself) and (2) this functor has special relations to itself, in that it is adjoint—in fact, *self-adjoint*. This perspective, developed formally in the following discussion, might even suggest, informally, that one think of the action of certain contravariant functors as a generalized sort of negation of structures. To develop these ideas, let us first recall some notions.

In general, regarding a poset as a category, recall that the colimit recovers the notion of supremum, while limit recovers that of infimum. In the case of the (co)limit over a diagram consisting of just two objects, we get the (co)product: the coproduct of objects x and y is the join (least upper bound) $x \vee y$, and the product is the meet (greatest lower bound) $x \wedge y$.

108. One position, in this context, might articulate the view that everything is what it is—as the individual thing it is—only on account of how it is *not* some other things, and accordingly try to take very seriously the idea that “all determination arises from a negation” (to paraphrase the old principle *omnis determinatio est negatio*, “every determination is a negation”). An opposing position might argue that negations always just describe *privations*, and a complete and accurate description of reality would not need to involve mention of any “negative entities.”

Definition 184 A poset \mathcal{P} is called a *join-semilattice* if each two-element subset $\{x, y\} \subseteq P$ has a join, denoted $x \vee y$. A poset \mathcal{P} is called a *meet-semilattice* if each two-element subset $\{x, y\} \subseteq P$ has a meet, denoted $x \wedge y$.

Definition 185 A poset is a *lattice* if it is both a join- and a meet-semilattice.

By induction, the binary case of joins and meets of a lattice can be extended so that every nonempty finite subset of a lattice has a join and a meet. In general, supposing we have \mathcal{P} a poset, then for all $x, y \in \mathcal{P}$,

$$x = \bigwedge \{x, y\} \iff x \leq y \iff y = \bigvee \{x, y\}.$$

For such a poset \mathcal{P} , for all $y \in P$, we can define the following:

Definition 186 (*A top as empty meet; bottom as empty join.*)

1. y is the empty meet, that is, $y = \bigwedge \emptyset$, iff it is a top (greatest) element; and
2. y is the empty join, that is, $y = \bigvee \emptyset$, iff it is a bottom (least) element.

Note that empty meets and joins need not exist. We could have two minimal elements, for instance, but no least element. A least element would have to be less than everything. An empty meet (top) is often written as \top or 1 while an empty join (bottom) is written as \perp or 0 .

If a lattice \mathcal{L} has a greatest element (top) 1 and a least element (bottom) 0 , which further satisfy that

$$0 \leq x \leq 1 \text{ for every } x \in \mathcal{L},$$

then we call the lattice a *bounded lattice*. Category-theoretically, this makes 0 and 1 the (unique) initial and terminal objects of the lattice, considered as a category. Thus, altogether, a lattice with 0 and 1 is a poset that, regarded as a category, has all finite limits and all finite colimits.

Posets are not guaranteed to have anything except the order \leq . Lattices, on the other hand, have all finite meets and joins. Further properties of lattices may obtain—for instance, distributive lattices obey an additional distributive law that brings them closer to logic.

Definition 187 A *distributive lattice* \mathcal{L} is a lattice in which the identity

$$x \wedge (y \vee z) = (x \wedge y) \vee (x \wedge z)$$

holds for all x, y, z . This identity gives the dual distributive law

$$x \vee (y \wedge z) = (x \vee y) \wedge (x \vee z).$$

We can also define the following notions:

Definition 188 For \mathcal{L} a bounded lattice, with least element 0 and greatest element 1 , we define a *complement* for an element x of the lattice as an element $a \in L$ such that $x \wedge a = 0$ and $x \vee a = 1$.

If a lattice is distributive, then a complement a , provided it exists, will be unique. In general, when it exists, the unique complement a of an element x is denoted by $a = \neg x$, meant to evoke a lattice-theoretic analogue of logical negation. This lets us define the following:

Definition 189 A *Boolean algebra* is a distributive lattice with 0 and 1 for which every element x has a complement $\neg x$.

Heyting algebras are examples of a still more general notion, namely of distributive lattices for which some members may lack complements. More explicitly,

Definition 190 A *Heyting algebra* H is a poset with all finite products and coproducts, and that is moreover Cartesian closed. Another way of describing such an H is as a distributive lattice with a least element 0 and a greatest element 1, expanded with an “implication” operation \Rightarrow , meaning that for any two elements p, q of the lattice, there exists an exponential q^p , usually written

$$p \Rightarrow q.$$

This operation is characterized by an adjunction, specifically

$$r \leq (p \Rightarrow q) \text{ iff } r \wedge p \leq q.$$

In other words, \Rightarrow is a binary operation on a lattice with a least element, such that for any two elements p, q of the lattice, $\max\{r \mid r \wedge p \leq q\}$ exists (where this latter set contains an element greater than or equal to every one of its elements, and such a least upper bound for all those elements r where $r \wedge p \leq q$ is what is denoted by $p \Rightarrow q$).¹⁰⁹

In the general case, in any Heyting algebra, we can define the “negation” of an element p as

$$\neg p := (p \Rightarrow 0).$$

Heyting algebras serve as models for intuitionistic propositional calculus—in which setting, variables are regarded as propositions, \wedge as “and,” \vee as “or,” and \Rightarrow as implication—and in that connection, it is sensible to think of “not p ” as effectively saying that “ p implies false.”

Note, moreover, how on account of the way \Rightarrow is defined, we can rewrite this as

$$q \leq \neg p \text{ iff } q \wedge p = 0,$$

revealing $\neg p$ to be the join of all those q whose meet with p in the lattice is 0, the least element. In any Heyting algebra H , we will have not just that

$$p \leq \neg \neg p,$$

but also that

$$p \leq q \text{ implies } \neg q \leq \neg p.$$

But this reveals how negation is just a contravariant functor from the Heyting algebra to itself! More explicitly,

Proposition 191 \neg is a functor $\neg : H \rightarrow H^{op}$ (and also $\neg : H^{op} \rightarrow H$). This functor is, moreover, adjoint to itself, since $p \leq \neg q$ iff $q \leq \neg p$.

Let us spell out this self-adjointness more explicitly. The first inequality, $p \leq \neg \neg p$, is immediate from $q \leq \neg p$ iff $q \wedge p = 0$, using the further fact that, for any Heyting algebra, $p \wedge \neg p =$

109. Another way to think of this $p \Rightarrow q$ is in the setting of the propositional calculus, where it is the weakest condition needed for the inference rule of modus ponens to hold, that is, to enforce that from $p \Rightarrow q$ and p we can infer q .

0 (this follows from the adjunctive definition of \Rightarrow together with the definition of negation as $\neg p = (p \Rightarrow 0)$). For the second, suppose $p \leq q$ in H . Then, $p \wedge \neg q \leq q \wedge \neg q$ and the right-hand side of this inequality is 0. So $p \wedge \neg q = \neg q \wedge p = 0$, and so by $q \leq \neg p$ iff $q \wedge p = 0$, we have that $\neg q \leq \neg p$. Incidentally, this moreover shows that

$$\neg p = \neg \neg \neg p$$

which we might call the “1 = 3” fact. This is a result of the contravariant functoriality of \neg and that $p \leq \neg \neg p$. For, suppose $p \leq \neg \neg p$. Then, by the contravariant functoriality inequality, we also have that $\neg \neg \neg p \leq \neg p$. And since $p \leq \neg \neg p$ holds for all p in H , it holds in particular for $\neg p$. Thus, $\neg p \leq \neg \neg \neg p$, giving the other side of the equality, and so, altogether, $1 = 3$.

That \neg , as a functor $H \rightarrow H^{op}$ and also $H^{op} \rightarrow H$, is adjoint to itself means that for all $p \in H$ and $q \in H^{op}(=H)$, we have the two-way rule

$$\frac{\neg p \leq^{op} q}{p \leq \neg q}$$

or

$$\frac{\neg p \geq q}{p \leq \neg q},$$

where the top (holding in H^{op}) holds if and only if the bottom (holding in H) does. For another way of seeing the truth of the fact that for all $x \in H$, $x \leq \neg \neg x$, notice that as an adjoint, letting $q = \neg p$, we must have

$$\frac{\neg p \geq \neg p}{p \leq \neg \neg p},$$

and since the top is always true, the bottom must be as well.

An adjoint is a kind of generalized inverse, and as such, an adjunction describes a kind of relaxation or weakening of the notion of equivalence. In the present situation, asking when it is in fact the case that $p = \neg \neg p$ (when the “do nothing” functor is *equal* to applying the negation functor twice) is like asking when the above adjunction happens to be an equivalence. If the relations \leq are replaced by $=$, then we get isomorphisms. This distinction captures the following well-known relation between Heyting algebras and Boolean algebras:

Proposition 192 A Heyting algebra H is Boolean (i.e., $\neg \neg x = x$ for all $x \in H$) if and only if the above adjunction is an equivalence.

This is a stricter requirement, and in general we need not have $x = \neg \neg x$. This requirement is something one might not always want to impose, and this is in fact one of the merits or utilities of working with the more general Heyting algebras. For any topological space X , the set $\mathcal{O}(X)$ of open sets of X forms a Heyting algebra; for instance, the opens in the real line accordingly form a Heyting algebra, but one that is not Boolean, since the complement of an open set is not necessarily open.

We could go on to dualize things and describe a dual notion, namely that of *co-Heyting algebras*, which support a corresponding but different notion of negation. The utility of considering such things can be motivated with another problem related to natural language. In many natural languages, for instance in English, one often has recourse to forms of negation that do not seem to be captured by, or behave as, the single negation operator of classical logic. Suppose someone is described to you as “not honest,” after they act in a particular way in a particular situation. This is not necessarily to say that they are “dishonest.” In natural language, we can deny that a person is honest in at least two distinct ways: (1) by asserting that someone is not honest (negating the predicable “to be honest”); or (2) by asserting that they are dishonest (negating the adjective “honest”). It is easy to appreciate, intuitively, how the second (“dishonest”) is a stronger form of negation than the first (“not honest”). Moreover, suppose your friend Abe is someone you would be willing to describe as “not dishonest.” This does not seem to convey the same thing as describing Abe as “honest.” Finally, while we expect it to be the case that Abe is either honest or not honest, it seems plausible to assert, as well, that he is neither honest nor dishonest. We would like to know how to capture these observations more formally, and describe formal relations between the different forms of negation, as applied to natural language. The story that follows begins to address this.

A little more generally, compare the sentence

It is false that not p ,

with the sentence

It is not false that p .

These sentences clearly do not say the same thing. The first indicates the *necessity* of p , while the second indicates its *possibility*. We can make sense of this in the context of a particular algebra called a bi-Heyting algebra, by interpreting the “it is false” in such sentences as the Heyting negation and the “not” as the corresponding “co-Heyting” negation. This setting will also allow us to define modal operators in terms of pairings of both negations.

Definition 193 A *co-Heyting algebra* is a poset whose dual is a Heyting algebra.

Unpacking this, we can equivalently observe that a co-Heyting algebra will be a bounded lattice expanded with a binary operation \searrow such that for every p, q, r , we have the adjunction rule

$$(p \searrow q) \leq r \text{ iff } p \leq q \vee r.$$

In other words, $p \searrow q = \bigwedge \{r \mid p \leq q \vee r\}$.

A corresponding unary negation operation \sim can then be defined by

$$\sim p := (1 \searrow p).$$

It follows that we have the following adjunction rule for this “negation” \sim :

$$\frac{\sim p \leq q}{1 = p \vee q}$$

Similar to how we have $p \wedge \neg p = 0$ in a Heyting algebra (though not necessarily $p \vee \neg p = 1$), notice how in a co-Heyting algebra we will have

$$p \vee \sim p = 1.$$

Likewise, similar to how the negation \neg is order-reversing and satisfies $x \leq \neg \neg x$ in a Heyting algebra, in a co-Heyting algebra the negation \sim is also order-reversing and it satisfies $\sim \sim x \leq x$. Moreover, just as in a Heyting algebra $p \vee \neg p$ is not necessarily the top (*true*), so in a co-Heyting algebra $p \wedge \sim p$ is not necessarily the bottom (*false*). In particular, then, in a co-Heyting algebra, we can thus define the generally nontrivial notion of the *boundary* of p , as

$$\partial p := p \wedge \sim p.$$

Incidentally, this recovers, in purely algebraic terms, the spatial and geometric notion of *boundary*, thus suggesting certain deep connections between logic and geometry and the study of space (connections explored in greater detail in the appendix). Bi-Heyting algebras—a bounded distributive lattice that is both a Heyting and a co-Heyting algebra—can accordingly be deployed to shed further light into some of these deep connections between logic and geometry.

We can work in the context of a bi-Heyting algebra and combine these negations to form our modal operators. In particular, we will let

$$\diamond p := \sim \neg p,$$

read as “possibly p .” This begins to suggest how we might formalize the fact that, for instance, John not being dishonest does not let us conclude, in general, that John is honest, but rather only that he is *possibly* honest. Similarly, we will let

$$\square p := \neg \sim p,$$

read as “necessarily p .”

Given such definitions, we can show that

Proposition 194 We have the adjunction $\diamond \dashv \square$.

Proof. To show that $\diamond \dashv \square$, we need only verify some equivalences, each of which more or less follows automatically from definitions. By definition, $\diamond = \sim \neg$ and $\square = \neg \sim$, so the adjunction just says that $\sim \neg \dashv \neg \sim$. But

$$\sim \neg p \leq q$$

just says that

$$1 \leq \neg p \vee q$$

which is of course equivalent to

$$1 \leq q \vee \neg p,$$

which itself is just to say that

$$\sim q \leq \neg p.$$

Using this last inequality, and the definition of \neg , we have

$$\sim q \wedge p \leq 0$$

or, equivalently,

$$p \wedge \sim q \leq 0.$$

This last line can be written

$$p \leq (\sim q \Rightarrow 0),$$

which is the same as saying that

$$p \leq \neg \sim q.$$

Altogether, this string of equivalences shows that

$$\sim \neg p \leq q \text{ iff } p \leq \neg \sim q,$$

which is precisely what is needed to show that $\sim \neg \dashv \neg \sim$ (or $\diamond \dashv \square$). □

Working in a bi-Heyting algebra and using the above definitions of \square and \diamond , we can moreover consider repeated applications of these operators. As a matter of simplifying computations involving repeated applications of the operators, we will define $\square_0 = \diamond_0 = \text{id}$, and

$$\square_{n+1} := \neg \sim \square_n, \quad \diamond_{n+1} := \sim \neg \diamond_n,$$

where \square_n is of course just the result of iterating (n times) the composition of $\neg \sim$, and \diamond_n by iterating (n times) the composition of $\sim \neg$. \square_n and \diamond_n are clearly both order-preserving, for all n , as is evident from the double fact that, in a Heyting algebra, \neg is order-reversing and satisfies $p \leq \neg \neg p$ for all p , and that, in a co-Heyting algebra, \sim is order-reversing and $\sim \sim p \leq p$. Moreover, we have

1. $\square_{n+1} \leq \square_n \leq \text{id} \leq \diamond_n \leq \diamond_{n+1}$ for all n ; and
2. $\diamond_n \dashv \square_n$ for all n .

Proof. 1. First, we know that, for any p in the bi-Heyting algebra, we have that $\neg p \leq \sim p$. From this, taking $p = \sim p$, we have that

$$\neg \sim p \leq \sim \sim p,$$

and since $\sim \sim p \leq p$, this gives that

$$\neg \sim p \leq p.$$

Moreover, $p \leq \neg \neg p$, and, using $\neg p \leq \sim p$ again, applied now to $p = \neg p$, we get that

$$\neg \neg p \leq \sim \neg p.$$

Altogether, this gives that

$$\neg \sim p \leq p \leq \sim \neg p.$$

By definition, then, this reads as

$$\square p \leq \text{id}(p) \leq \diamond p,$$

and further iterating this, letting $p = \square_n p$, and then $\diamond_n p$, gives the main result. □

Proof. 2. We want that $\diamond_n \dashv \square_n$ for all n . But since adjoints compose, this follows from the preceding result, by iterating. □

In a bi-Heyting algebra where countable suprema and infima exist satisfying

$$\frac{b \leq \bigwedge_n a_n}{b \leq a_n \text{ for all } n}$$

and

$$\frac{\bigvee_n a_n \leq b}{a_n \leq b \text{ for all } n},$$

we can further define

Definition 195

$$\begin{aligned} \Box p &:= \bigwedge_n \Box_n p, \\ \Diamond p &:= \bigvee_n \Diamond_n p. \end{aligned}$$

Then, for a bi-Heyting algebra that has countable suprema and infima satisfying the above rules, the modal operators \Box and \Diamond are such that $\Box p$ will be the largest complemented x such that $x \leq p$, and $\Diamond p$ will be the smallest complemented x such that $p \leq x$. This moreover realizes the fact that \Box and \Diamond are both order-preserving and that $\Box p \leq p \leq \Diamond p$.

More generally, for any bounded distributive lattice and operators \Box and \Diamond defined just as above, the same properties will hold of \Box and \Diamond . In particular, it can be shown that $\Diamond \dashv \Box$, and thus

1. $\Box \leq \text{id} \leq \Diamond$,
2. $\Box \Box = \Box, \Diamond \Diamond = \Diamond$,
3. $\text{id} \leq \Box \Diamond$,
4. $\Diamond \Box \leq \text{id}$,
5. $\Diamond(\phi \wedge \Box \psi) = \Diamond \phi \wedge \Box \psi$.

7.2.3 Adjoint Modalities in the Qua Category

Example 196 For another realization of some of the ideas explored in the previous section, recall the *interpretation* functor on the *qua* category **Qua** from example 38 (chapter 2).¹¹⁰ Applying all this to the particular subcategory **A** from earlier, an interpretation amounts to a set X together with a set of predicates of X , where by “predicate of X ” is just meant a family $\{\phi_A\}_{A \in \text{Ob}(\mathbf{A})}$ of subsets of X with the functorial property

$$\text{if } x \in \phi_A \text{ (i.e., } x \in_A \phi) \text{ and } A' \rightarrow A \in \mathbf{A}, \text{ then } x \in \phi_{A'}.$$

And since **A** was just the comma or slice category $\boxed{\text{a scf}} \downarrow \mathbf{CN}$, so that all aspects are of the form $\boxed{\text{a scf}} \text{ qua } \boxed{\mathbf{B}}$, an interpretation of **A** just associates to every aspect the same set X . Notice that we can collect together such families of predicates (as subfunctors of X) into the set $\mathcal{P}(X)$ of all predicates of X . These predicates in fact form a bounded distributive lattice with two negations—in fact, a bi-Heyting algebra—as

$$(\mathcal{P}(X), \leq, \vee, \wedge, 1, 0, \neg, \sim),$$

110. Again, this material on the *qua* category is ultimately derived from Reyes et al. (1999); the reader who desires to pursue these matters further than the discussion here can find many more interesting details in that paper.

where there is the natural ordering

$$\phi \leq \psi \text{ iff } \forall A \in \mathbf{A}, \forall x \in X, x \in_A \phi \Rightarrow x \in_A \psi.$$

Also, as expected, we have

$$x \in_A (\phi \vee \psi) \text{ iff } x \in_A \phi \text{ or } x \in_A \psi$$

and

$$x \in_A (\phi \wedge \psi) \text{ iff } x \in_A \phi \text{ and } x \in_A \psi.$$

0 (or \perp) is the predicate “false,” the bottom element of the order, while 1 (or \top) is the predicate “true,” the top element of the order. Given a predicate ϕ , we define two negations, $\neg\phi$ and $\sim\phi$, which are in fact two new predicates (i.e., have the requisite functorial property):

$$x \in_A \neg\phi \text{ iff } \forall A' \rightarrow A \in \mathbf{A} \ x \notin_{A'} \phi$$

and

$$x \in_A \sim\phi \text{ iff } \exists A \rightarrow A' \in \mathbf{A} \ x \notin_{A'} \phi.$$

Applied to the predicate “honest,” for instance, we have the natural reading: “ \neg honest” as “dishonest,” and “ \sim honest” as “not honest.”

As we saw in our earlier discussion of these same negations, we have the following adjunctions for our negations (which hold for arbitrary properties ϕ, ψ):

$$\begin{array}{c} \psi \leq \neg\phi \\ \psi \wedge \phi = 0, \\ \sim\phi \leq \psi \\ 1 = \psi \vee \phi. \end{array}$$

Notice also that for every property ϕ , it is a consequence of the above that $\phi \wedge \neg\phi = 0$ and $\phi \vee \sim\phi = 1$. However, $\phi \wedge \sim\phi$ need not be 0, and similarly, $\phi \vee \neg\phi$ is not necessarily 1.

Applying all this, suppose we return to modeling a discussion and decision regarding your friend Abe’s honesty. We can assume the aspects considered relevant to Abe’s honesty have been agreed upon, and likewise, agreement has been achieved on his honesty under each of the relevant aspects, that is, for every aspect $A \in \mathbf{A}$, we assume it is known whether $Abe \in_A \text{ honest}$ or $Abe \notin_A \text{ honest}$. If Abe fails to be honest under every subaspect of one of the aspects, say $F = \text{family man}$, then we would say that “Abe is dishonest” under that aspect, that is, *qua* family man. By contrast, we would say that “Abe is not honest” under the aspect F precisely when he fails to be honest with respect to one of the superaspects of F . The ultimate judgment regarding Abe’s honesty is then obtained by restricting to the global aspect, G (or “*qua* scf”), where this means that Abe is “honest,” “not honest,” or “dishonest” precisely when $Abe \in_G \text{ honest}$, $Abe \in_G \sim \text{honest}$, $Abe \in_G \neg \text{honest}$, respectively. In more detail,

$Abe \in_G \text{ honest}$ iff $\forall A \ Abe \in_A \text{ honest}$, that is, “Abe is honest iff Abe is honest under any aspect.”

$Abe \in_G \sim \text{honest}$ iff $\exists A \ Abe \notin_A \text{ honest}$, that is, “Abe is not honest iff Abe fails to be honest under at least one of the aspects.”

$Abe \in_G \neg \text{honest}$ iff $\forall A \ Abe \notin_A \text{ honest}$, that is, “Abe is dishonest iff Abe fails to be honest under every one of the aspects.”

If Abe is honest, then he cannot be dishonest (and conversely), that is, $\phi \wedge \neg\phi = 0$. However, that $\phi \vee \neg\phi$ is not necessarily 1 means, of course, that it is not always the case that Abe is either honest or dishonest. One scenario in which this might occur would be where Abe is honest under the aspect S (Abe is honest *qua* student) but fails to be honest in all the other aspects. This reveals how the negation is not, in general, Boolean. Yet we can note that \sim is Boolean globally. This means that for the “global aspect,” we will have that Abe is either honest or not honest, but not both, that is, $honest \vee \sim honest = 1$ and $honest \wedge \sim honest = 0$. However, it may occur that Abe is both honest and not honest under the very same aspect (as long as this is not the global aspect).

Thus, notice that even though Abe is “not honest” under a particular aspect precisely when he is not honest under every aspect, for all aspects other than the global one, there is a difference between “not honest” under an aspect and *failing* to be honest under that aspect, for the former has the functoriality property: if Abe is not honest under aspect A , then he is not honest under any subspect $A' \rightarrow A$. Failing to be honest under a given aspect, by contrast, is simply the absence of Abe’s honesty under that aspect. Absence of honesty under an aspect is not functorial, so there may be an absence of Abe’s honesty under an aspect and the presence of Abe’s honesty under another.

Suppose Abe is not dishonest under an aspect. What can we conclude from this? Not that Abe is honest; rather, only that he is *possibly* honest. Abe is not dishonest under a given aspect precisely when he is honest under at least one aspect, as the following string of equivalences reveal:

$$\frac{\frac{\frac{Abe \in_A \sim \neg honest}{\exists A \rightarrow A' Abe \notin_{A'} \neg honest}}{\exists A \rightarrow A' \exists A'' \rightarrow A' Abe \in_{A''} honest}}{\exists A' Abe \in_{A'} honest}$$

Restricting to the global level $A = G$, gives

$$\frac{Abe \in \sim \neg h}{\exists A' Abe \in_{A'} h},$$

or “Abe is not dishonest iff Abe is honest under at least one aspect.” Since this works for any predicate, this suggests we define our modal operator

$$\diamond \phi := \sim \neg \phi,$$

read as “possibly ϕ .”

Similarly, we can calculate

$$\frac{\frac{\frac{Abe \in_A \neg \sim honest}{\forall A' \rightarrow A Abe \notin_{A'} \sim honest}}{\forall A' \rightarrow A \exists A'' \rightarrow A' Abe \in_{A''} honest}}{\forall A' Abe \in_{A'} honest}.$$

Again, letting $A = G$ the global aspect, this becomes

$$\frac{Abe \in \neg \sim honest}{Abe \in honest},$$

that is, Abe is necessarily honest under a given aspect iff he is honest under every aspect. In other words, $x \in_A \square \phi$ iff for all aspects A' , $x \in_{A'} \phi$ iff $x \in_G \phi$. This suggests we define a further modal operator,

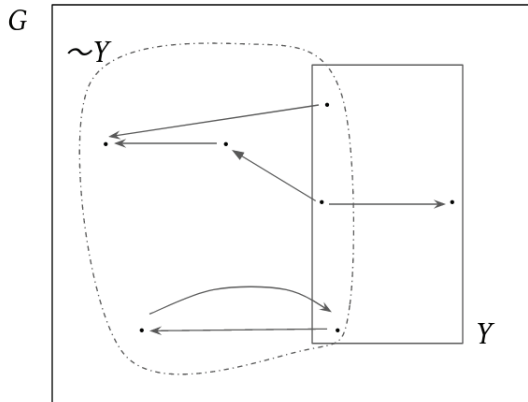
$$\square \phi := \neg \sim \phi,$$

read as “necessarily ϕ .” Observe that $\diamond\phi$ and $\Box\phi$, defined thus, are clearly themselves predicates of X . Moreover, they are themselves adjoint, $\diamond \dashv \Box$, and as such satisfy the (in)equalities we isolated at the end of section 7.2.2.

7.2.4 Adjoint Modalities in Graphs

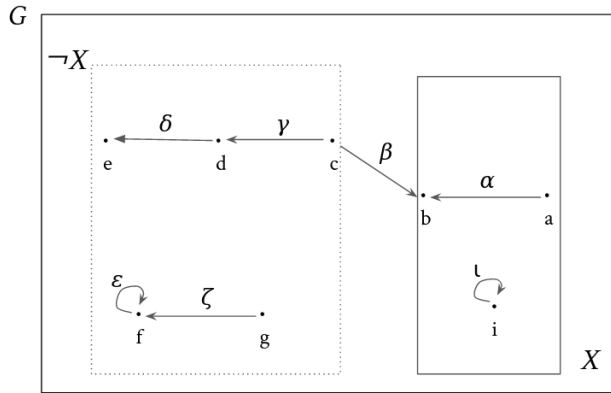
Example 197 The modalities discussed in the last few sections are particularly well illustrated and tangible in the context of graphs and their subgraphs. If we have a directed multigraph G , the lattice of subgraphs of G constitutes a bi-Heyting algebra,¹¹¹ where a subgraph X of G is a directed multigraph (i.e., consisting of a subset X_0 of the vertices of G and a subset X_1 of the edges of G) such that every edge in X_1 has both its source and target in the vertex set X_0 . It is clear that we can take unions and intersections of subgraphs, but what it means to take “complements” is not as evident. The set-theoretical complement X^c of a subgraph X will not work, since in general it will not even be a graph, as we might wind up with edges whose source or target is missing in the set X^c . There are two obvious ways to address the insufficiencies with this complement operation, though: we can either discard such problem edges or, on the other hand, we can retain them and “complete” them by adding their sources and targets in the underlying graph. The first option in fact gives rise to our Heyting negation $\neg X$, and the second to the co-Heyting $\sim X$.

It is instructive to see these notions “at work” and to make explicit computations with them. So observe that given the graph G together with its subgraph Y , as displayed below, then for $\sim Y$ we will have

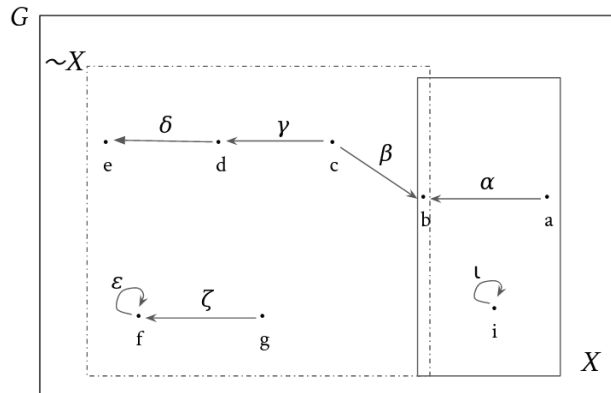


Running this again, we compute $\sim\sim Y$ (which is, importantly, not identical to $Y!$),

111. This is actually a special case of something that will be discussed further in chapter 10, namely that any presheaf topos is bi-Heyting. We know that the category of (multi)graphs can be represented as $\mathbf{Set}^{\mathbf{C}^{op}}$ for \mathbf{C} the index category of two objects and two nontrivial morphisms between those objects; moreover, it can be shown that the lattice of “subobjects” of any object of a bi-Heyting topos is itself a (complete) bi-Heyting algebra.

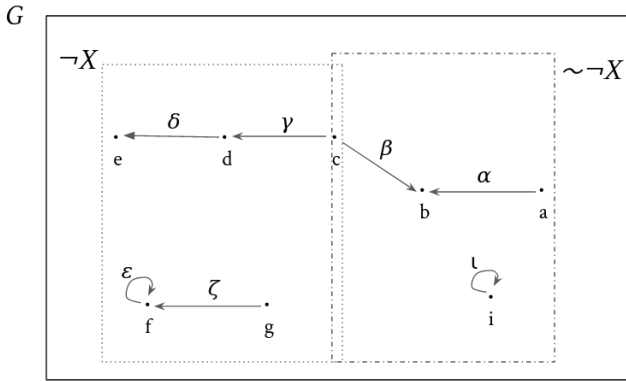


$\sim X$, for its part, is here different from $\neg X$, and yields the smallest subgraph whose union with X gives all of G :

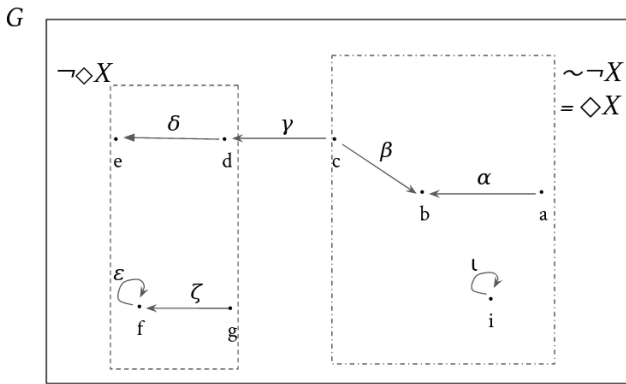


Incidentally, notice that the boundary of X , $\partial X = \sim X \wedge X$, is not the empty subgraph (which functions as 0), but is the sole vertex b . Intuitively, it makes sense to think of the vertex b as the “boundary” of X , since it liaisons between the “inside” of X and the “outside,” as there is an arrow, namely β coming in to X from the outside of X via the target b .

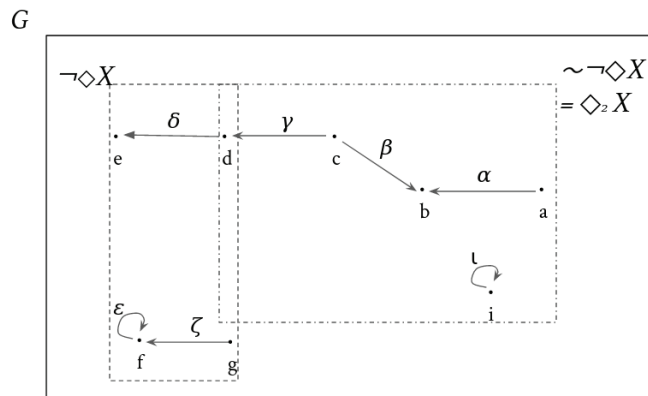
Let us now compute $\diamond_1 X$, that is, $\sim \neg X$:



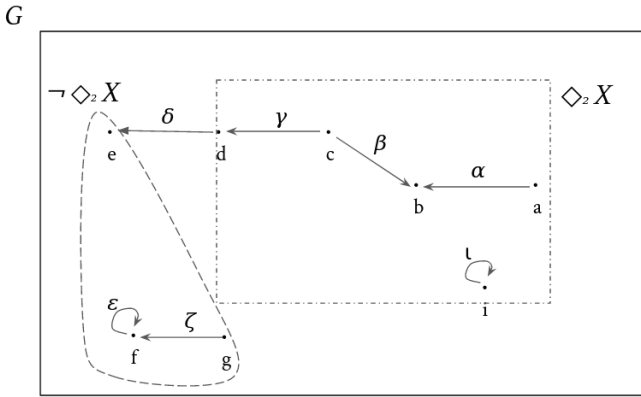
To compute $\diamond_2 X$, we must first compute $\neg \diamond_1 X$



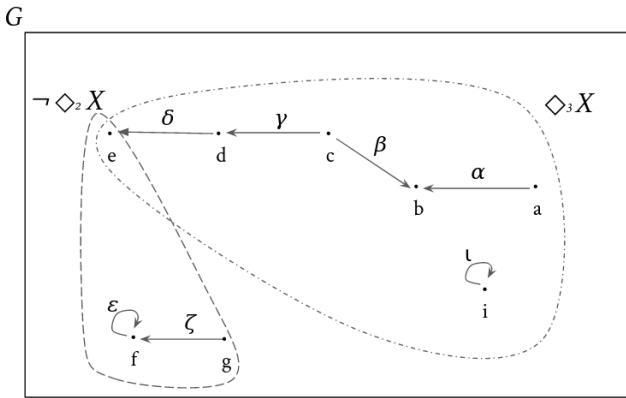
and then $\sim \neg \diamond_1 X$, that is, $\diamond_2 X$,



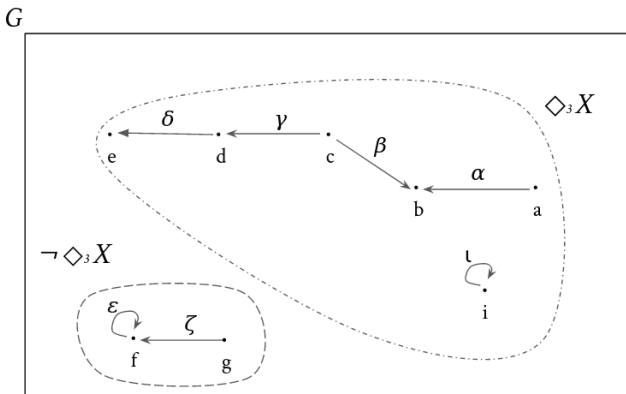
Running this one more time, we first get $\neg \diamond_2 X$,



and then \sim of this, that is, $\sim \neg \diamond_2 X = \diamond_3 X$,

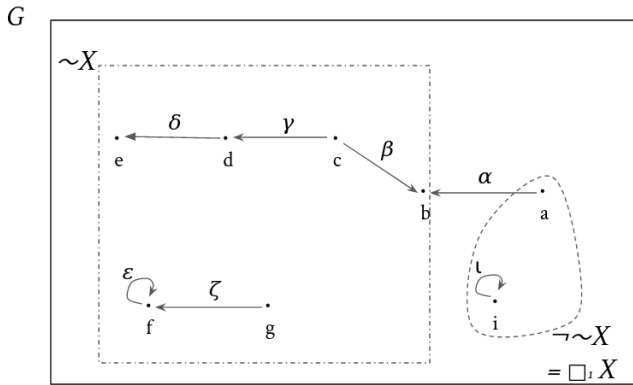


At this point, something interesting happens. If we run \neg on $\diamond_3 X$, we get

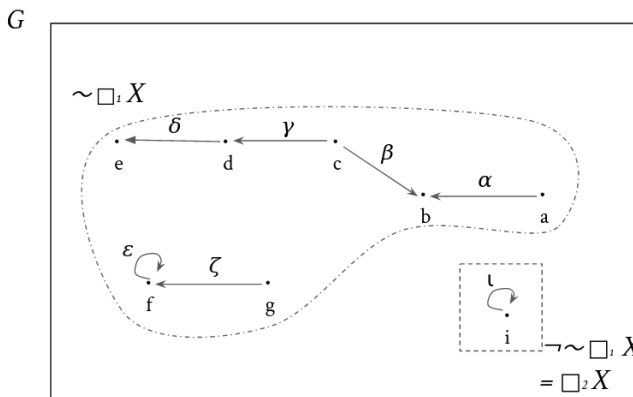


and then \sim of the above, which gives us $\diamond_4 X$, is revealed to be the same as $\diamond_3 X$. Thus, this iterative operation stabilizes at $\diamond_3 = \diamond_4$, and we have that $\diamond X = \diamond_3 X$, having captured, in a subgraph, all those elements of the graph G that can be reached from X through some path, which is clearly something of a picture of the *possibility* of X (or, perhaps more accurately stated, of what is *possible for* X). This illustrates a more general feature as well, namely that every arrow or vertex that is connected to X via some path will end up in $\diamond X$ after a finite number of steps. In general, applying the operator \diamond_n to the subgraph should be thought of as capturing those elements connected with X within n paths. Notice also that taking \neg of \diamond_3 is the same as taking \sim of \diamond_3 . As such, $\diamond X$ has no boundary (and no edge going out of it), as $\partial \diamond X = \diamond X \wedge \sim \diamond X$ is the empty subgraph. In general, in the land of graphs, taking the boundary of a subgraph X yields the subgraph whose elements are connected to the outside of X .

We can perform similar computations, in reverse order, to compute $\square X$, which will supply us with a subgraph whose elements are those that are *not* connected to the outside. First, we compute \neg of $\sim X$, to get $\neg \sim X = \square_1 X$,



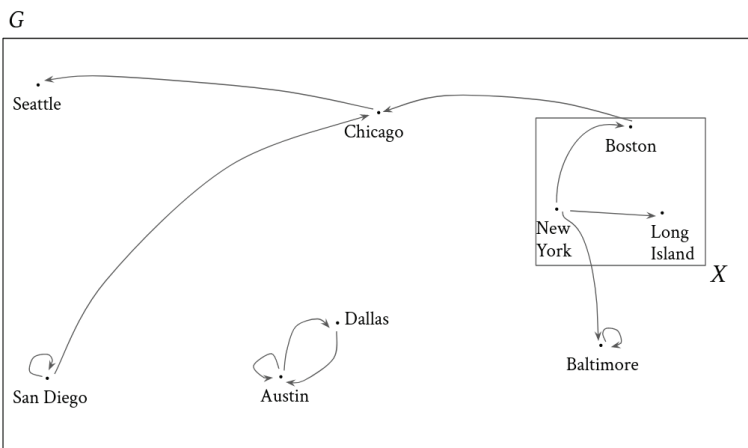
Then, taking \sim of this, we get $\sim \square_1 X$, and stripping away one more layer, by taking \neg of the result, we end up with $\neg \sim \square_1 X = \square_2 X$:



and we have stability, as any further iterations $\square_{2+n}X$ will just reduce back to \square_2X . Notice that what is left, $\square X$, just consists of those elements of X that are not connected to the outside (of X in G), which seems to align with some intuitions we have about the notion of “necessity” for X .

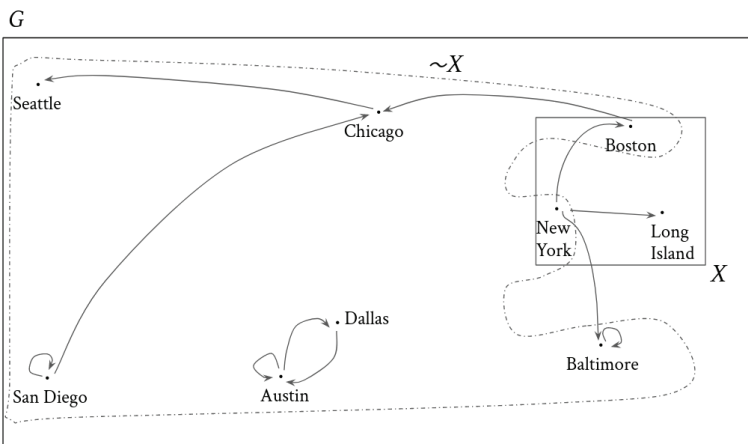
Altogether then, and for a general graph, $\diamond X$ supplies the elements of the ambient graph G that can be reached from X via some path, while $\square X$ has those elements in X not connected, via any path, to the outside. Note, finally, that both the subgraphs $\diamond X$ and $\square X$ are complemented sums of connected components.

Exercise 20 Consider the following graph G of routes, with subgraph X corresponding to some region of the northeast (including the nodes Boston, New York, and Long Island, together with the indicated routes between them):

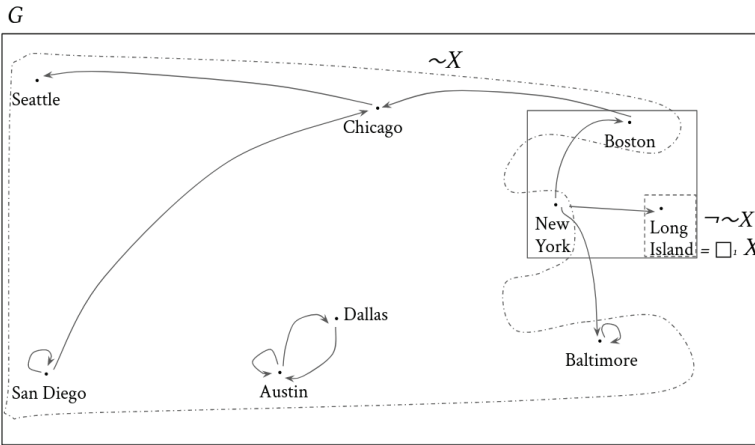


Compute $\diamond X$, $\square X$, and ∂X . Then consider, via this example, how the boundary operator ∂ interacts with the modal operators.

Solution First, notice that

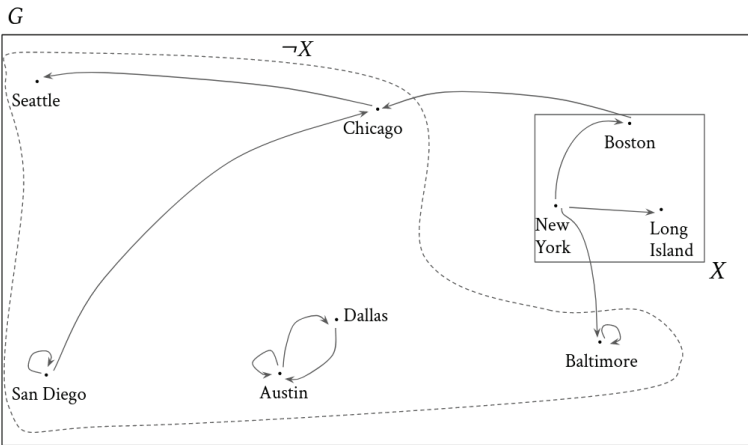


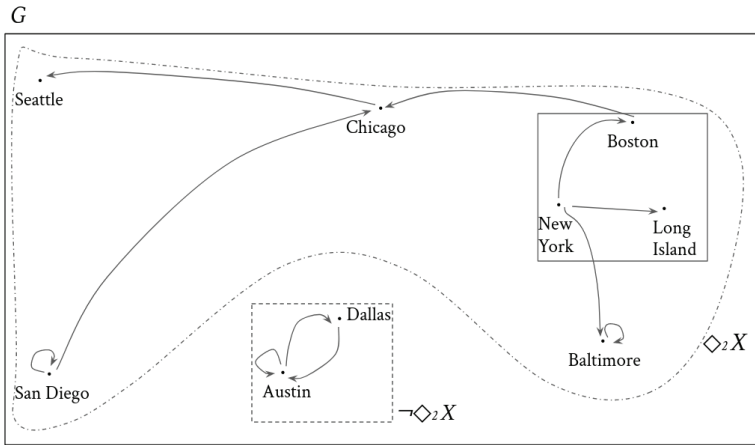
and



Things stabilize here, with $\square_1 = \square$. $\square X$ includes those parts of X that have no connection to the outside of X . The meaning of $\square X = \{ \text{Long Island} \}$ is this: if one is in X and ends up in Long Island, one will never get out of X —having arrived there, one is *necessarily* in X .

What about $\diamond X$, or “possibly” X ?





Taking \sim of $\neg\Diamond_2X$ just returns \Diamond_2X , so we achieve stability at $\Diamond_3 = \Diamond_2$, and so \Diamond_2X gives us $\Diamond X$, a picture of the “possibility” of X . Intuitively, this makes sense, as $\Diamond X$ supplies those parts of G that can be connected, directly or indirectly, with some part of X . For instance, then, even though San Diego is not directly reachable from any city in X or any city reachable by X , anyone in San Diego can get to Chicago, and Chicago *is* reachable from X . In other words, a person from X might meet someone from San Diego in Chicago or Seattle—and so, for someone from X , San Diego may form part of their picture of reality. Dallas and Austin, by contrast, are inaccessible to X —given the graph above, someone from anywhere in X could never meet anyone from Dallas or Austin.

Finally, consider $\partial X = \sim X \wedge X$. As one can see, this will give the “vertices” Boston and New York. Intuitively, this makes sense, as these cities are those parts of X that mediate between the “inside” of X (as parts of X) and the “outside” in G . Seeing how the boundary operator ∂ interacts with the modal operators can further solidify the intuitiveness of the reading of \Diamond as “possibility” and \Box as “necessity,” even in contexts like that of graphs. As one can easily verify,

$$\partial\Box X = \sim\Box X \wedge \Box X = \Box X$$

which confirms the intuition that the boundary of what is necessarily X —that is, those parts of X that have no connection to the outside of X —is just trivially $\Box X$ itself. Also,

$$\partial\Diamond X = \sim\Diamond X \wedge \Diamond X = \emptyset$$

or, more accurately, the empty subgraph. Intuitively, this realizes the idea that the “world” of what is *not possible* for X has empty overlap with what’s *possible* for X .

7.2.5 Adjoint Modalities in Topology

Example 198 Recall from chapter 4 (and see the appendix for) the extended discussion of open and closed sets of a topology, and the associated operations of taking the interior and closure. If we let $\mathcal{O}(X)$ denote the open sets, and $\mathcal{C}(X)$ the closed sets, of some space X , then **int** can be regarded as an inclusion-preserving map from $\mathcal{C}(X)$ to $\mathcal{O}(X)$, and **cl** as an inclusion-preserving map from $\mathcal{O}(X)$ to $\mathcal{C}(X)$. Building on the work of chapter 4 and the appendix, we can leave it to the reader to verify that these maps satisfy that, for any U

open and A closed,

$$\mathbf{cl}(U) \subseteq A$$

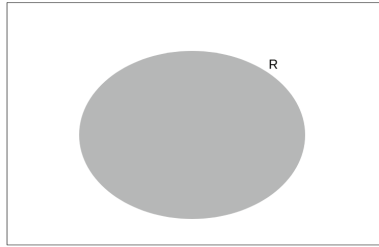
iff

$$U \subseteq \mathbf{int}(A).$$

This situation in fact just describes, in category theoretic terms, that \mathbf{cl} is left adjoint to \mathbf{int}

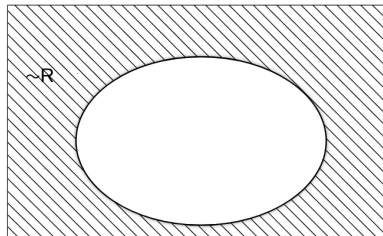
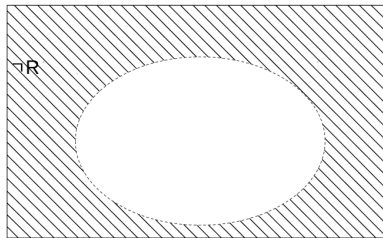
$$\mathcal{O}(X) \begin{array}{c} \xrightarrow{\mathbf{cl}} \\ \perp \\ \xleftarrow{\mathbf{int}} \end{array} \mathcal{C}(X).$$

Moreover, suppose given the region R



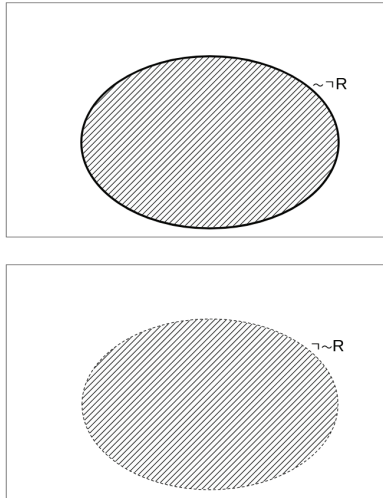
There are two ways of negating (taking the complement of) a part R of a space, where we may assume we initially know nothing of whether or not R includes its boundary. Namely:

1. do not include the boundary, that is, take $\neg R$;
2. do include the boundary, that is, take $\sim R$.



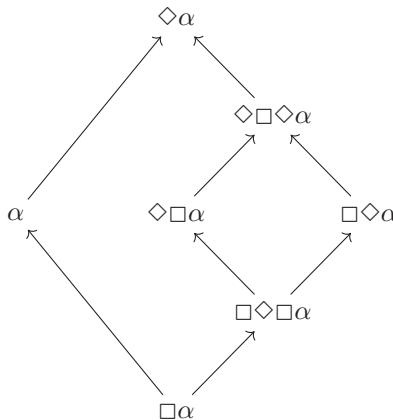
The second operation acts as the closure of the complement. By carrying out these operations again on the results, we can in particular get

1. $\sim \neg R$ (closure of R), or $\diamond R$;
2. $\neg \sim R$ (interior of R), or $\square R$.



This starts to hint at some of the profound connections between the modal operators \Box and \Diamond , on the one hand, and the topological operators **int** and **cl**, on the other—connections explored in the appendix. The more general story in terms of the adjointness of both pairs further allows us to exhibit and unify otherwise seemingly arbitrary results from each special area. The following displays one such case.

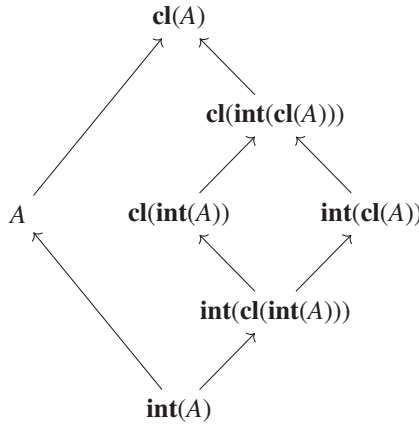
In modal logic, a *modality* is any sequence of zero or more monadic operators that involves \neg , \Box , and \Diamond . Iterated modalities include more than one such operator. One is sometimes asked to prove that the logic **S4** has (up to equivalence) exactly fourteen modalities, specifically: (1) id; (2) \Box ; (3) \Diamond ; (4) $\Box\Diamond$; (5) $\Diamond\Box$; (6) $\Box\Diamond\Box$; (7) $\Diamond\Box\Diamond$; together with the negations of each (i.e., adding \neg to each), making for fourteen in total. In the affirmative case (unnegated), the seven modalities are related by implication as displayed in the following diagram (where the arrow should be read “implies”):



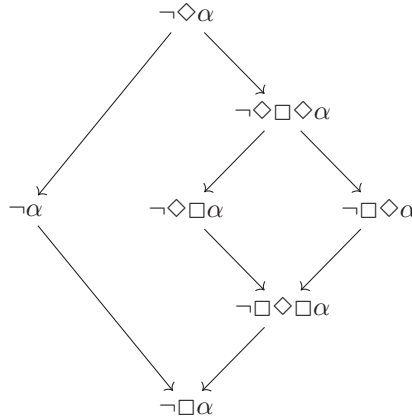
In the context of topology, one can verify that

$$\mathbf{cl}^2 = \mathbf{cl}, ((-)^c)^c = \text{id}, \mathbf{int} = (\mathbf{cl}((-)^c))^c, \mathbf{int}^2 = \mathbf{int}, \mathbf{int}((-)^c) = (\mathbf{cl}((-)^c))^c, \mathbf{cl}((-)^c) = (\mathbf{int}((-)^c))^c,$$

and there is the following diagram relating the closure and interior operators, and combinations thereof,



Taking complements gives a dual diagram. Similarly, by negating each of the modalities in the first diagram of modalities, and accordingly reversing the order of implication, we can display the other seven modalities and their mutual relations:



This can be compared to the diagram governing the topological operators, where \neg becomes complement $(-)^c$, and \square is replaced with **int** and \diamond with **cl**.

In topology, there is a somewhat mysterious result—called Kuratowski’s 14-set theorem—which says that in a topological space, fourteen is the maximum possible number of distinct sets that can be generated from a fixed set by taking closures, interiors, and complements (or just by taking interiors and complements, or closures and complements). The above treatment of these matters, together with the connections explored in the appendix, gives a unified framework for understanding why this, and the fourteen modalities results, should be true, and begins to suggest how they form part of a single story. As it turns out, while this number fourteen is generally a *maximum*, a topological space X will

actually contain fourteen such sets when it is “sufficiently rich,” in particular when it contains a copy of the Euclidean line—which fact has close connections to the fact (explored in the appendix) that the logic **S4** models the Euclidean line.¹¹³

7.3 Some Additional Adjunctions and Final Thoughts

More examples of important adjunctions will arise organically throughout the book, and the facts learned in this chapter will be put to use. At this point, there are a number of more advanced category-theoretic results and constructions that we could pursue. But with a good working understanding of categories, functors, natural transformations, and adjunctions—built on a number of carefully selected examples and applications—the reader is already well equipped to delve deeper into sheaves.

For now, this chapter ends by sketching a few more examples, the last two of which are left deliberately somewhat vague, while being correct “in spirit” (and, in fact, they can be developed to be formally correct). These are meant to provide what I hope are some engaging examples of adjunctions, while encouraging more fastidious readers to work out the unspecified details on their own. The chapter ends with a brief Philosophical Pass.

Example 199 Recall how every (directed) graph gives rise to a category and how every category is itself a directed graph with some extra data concerning composition. There is in fact a “free-forgetful” adjunction $F \dashv U$ here

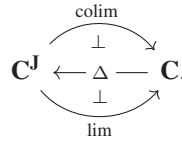
$$\mathbf{Cat} \begin{array}{c} \xrightarrow{U} \\ \xleftarrow{F} \end{array} \mathbf{Grph}.$$

The functor F gives the free category on a graph first mentioned in example 34 (chapter 2), which takes vertices for its objects, edges for its morphisms, adds identity arrows at each vertex, and lets the set of morphisms between two vertices be the set of all finite paths between the two. The functor U , for its part, just assigns to a category its underlying graph—it does this by forgetting all the data and structure of the category except for the objects (sent to vertices) and morphisms (sent to edges), leaving out any information or conditions pertaining to identities and composition, as these things are not native to graphs.

Example 200 Universal concepts like colimits and limits, initial objects and terminal objects, can be phrased entirely in terms of adjoint functors. As mentioned already, one of the advantages of this adjunction perspective is that the (co)limit of *every* **J**-shaped diagram in **C** can be defined all at once, rather than just taking the (co)limit of a particular **J**-shaped diagram $X : \mathbf{J} \rightarrow \mathbf{C}$.

Recall the notion of diagrams of shape **J** and the constant diagram functor $\Delta : \mathbf{C} \rightarrow \mathbf{C}^{\mathbf{J}}$, which sends each object $c \in \mathbf{C}$ to the associated constant functor, defined as the composite functor $c \circ t : \mathbf{J} \rightarrow \mathbf{C}$, where $t : \mathbf{J} \rightarrow \mathbf{1}$ is the unique functor to the terminal category and c is the object represented by the functor $c : \mathbf{1} \rightarrow \mathbf{C}$. Then, given an object X of $\mathbf{C}^{\mathbf{J}}$, we want an object of **C**. Indeed, it turns out that a category **C** admits all limits of diagrams indexed by a small category **J** iff the constant diagram functor Δ admits a right adjoint, and that **C** admits all colimits of diagrams indexed by **J** iff Δ admits a left adjoint:

113. For more on Kuratowski’s 14-set theorem, see Sherman (2004). For more on the connections between the notions of modal logic and features of topology, see the appendix.



Example 201 There is a connection between how the world *appears* to an agent and what that agent *believes* to hold of their world. But “appears” and “believes” are not quite inverses of one another. Instead, we might conjecture that

$$\text{appearance} \dashv \text{belief},$$

in the sense that there is an adjunction (in a slogan, “belief as the right adjoint of appearance”)

$$\frac{A_\alpha(m) \leq m'}{m \leq B_\alpha(m')}$$

realizing, effectively, how all that appears to an agent to hold at, or given, state m entails state m' if and only if whenever m holds in the “real world,” this entails that all that the agent *believes* to hold on the assumption that m' holds does in fact hold. Thus, in general, $B_\alpha(m)$ would stand for agent α ’s belief at m and will consist of those propositions that agent α *believes to hold* whenever m holds.

Example 202 Both small-scale and large-scale projects, such as those in research or development, require resources. Resource allocation (through grants, investment funding, contracts, etc.) requires a detailed plan (for how those resources are to be spent), especially as the project increases in scale. If **Rsrc** is a category consisting of relevant resources, so that objects are resources (e.g., for simplicity, different-sized checks or bags of money) and morphisms are given by a natural relation between those resources (e.g., \leq in the case of a uniform money-valuation of the different resource objects), and if **ProjPlan** is a category consisting of project tasks, given some natural ordering (e.g., by order of priority in the carrying out of the plan), then we might consider the functor

$$V : \mathbf{Rsrc} \rightarrow \mathbf{ProjPlan}$$

that maps a resource r to the collection of plans p_i that are *viable* given that resource, and the functor

$$N : \mathbf{ProjPlan} \rightarrow \mathbf{Rsrc}$$

taking a project task p to all those resources that are necessary to complete the task (which, depending on how **Rsrc** is structured, say in a simple case of “costs,” might just amount to returning an interval bounded by the *least* cost for which the task could be carried out, and including all other more ample amounts).

We would probably not expect V and N to construct strict *inverses* to one another, for we do not expect that, for any given resource r , a list of necessary resources for those plans that are deemed viable given r would be *equal* to r . Though we might expect that, among the resources, $r \leq NV(r)$. Similarly, we would not expect that, for a given project task p , the result of applying N to p and then V to $N(p)$, would always *equal* the same task p . Yet we would expect $VN(p) \leq p$ in **ProjPlan**. This suggests that we have an adjunction,

$$\mathbf{Rsrc} \begin{array}{c} \xrightarrow{V} \\ \perp \\ \xleftarrow{N} \end{array} \mathbf{ProjPlan}.$$

At this point, having seen a variety of examples and explored some of the basic facts of the general theory of adjoint functors (e.g., adjoints are unique up to unique natural isomorphism, adjunctions can be composed, and so on), readers should feel more comfortable with the idea of adjunctions. This simple concept—fundamentally consisting of an opposing pair of functors with a special relationship to one another—is incredibly powerful and useful. Adjoint functors really are ubiquitous throughout mathematics,¹¹⁴ and the examples considered in this chapter are but a very small fraction of panoply of adjunctions that interest mathematicians. Adjoint functors generally tell us very important things. For instance, if a given functor has a left adjoint, then (being a right adjoint of the adjoint pair) it will commute with limits (i.e., it is continuous); if it has a right adjoint, then (being a left adjoint of the pair) it commutes with colimits (i.e., it is cocontinuous). Moreover, we know that the left (or right) adjoint to a functor, provided it exists, will be unique up to natural isomorphism. But adjoint functors need not exist. Given a functor, it need not admit an adjoint. Readers who may very well appreciate the utility and interest of adjoint functors might still be wondering: how do we find such things and how can we know if they even exist? In searching for adjoint pairs, it can often be useful to keep in mind those general situations where adjoints do exist—for instance, how when dealing with categories that consist of algebras (groups, vectors, etc.), forgetful functors will have left adjoints—and consider if they apply. More abstractly, we can *use* some of the general results to help us with the question. In particular, if your given functor T can be shown not to preserve either limits or colimits (i.e., not to be continuous or cocontinuous), then it cannot be either the left or right adjoint of an adjoint functor pair. In other words, (co)continuity gives us a necessary condition for a functor to have a left (right) adjoint. What about sufficient conditions? As it turns out, it can be shown that a functor $T : \mathbf{A} \rightarrow \mathbf{S}$ admits a left adjoint iff for each $s \in \mathbf{S}$ the comma category $s \downarrow T$ has an initial object, which effectively means that the problem of finding a left adjoint to a continuous functor can be treated as a problem of finding an initial object in the associated comma category given for each object of \mathbf{S} .¹¹⁵ There are additional conditions—supplied by the so-called “adjoint functor theorems”—that will apply in certain contexts and give sufficient conditions for an adjoint to exist.

In later chapters—especially chapter 11—we will explore important adjunctions in greater depth, in the course of which we will see how, in situations where we can interpret constructions in terms of presheaf (or sheaf) categories, there will automatically be certain pairs of adjoint functors.

7.4 Philosophical Pass: The Idea of Adjointness

Box 7.1

The Idea of Adjointness

114. See Mac Lane’s oft-cited slogan, “Adjoint functors arise everywhere” (Mac Lane 1998).

115. For a proof and discussion of this, as well as the other adjoint functor theorems referred to in the following sentence, see Riehl (2016, Lemma 4.6.1 and Epilogue).

Adjoint functors are perhaps the most decisive concept in category theory. In addition to being found all over mathematics, adjoint functors frequently arise from constructions that have a certain universal property, where the constituent functors of an adjoint pair give the “most efficient” solution to a problem. To put it another way: many constructions with universal properties can be translated into statements of adjointness.

Beyond the connections to universal properties, adjunctions exhibit what category theory is really all about. To somewhat overstate the point: categories are not what is important. And the purpose of category theory is not to *categorize*, if that means what it usually means—creating an inventory of sorts that we can consult to ensure that things are put where they belong. What *is* important—the reason why we care about categories in nearly all cases—is supplied by the functors between them and the relations between those functors. Many mathematical structures will have all sorts of important relations to other structures, relations that cannot be reduced to questions of whether the underlying structures they are relating are the “same” (in the sense of isomorphism or equivalence). As just the right formal substitute for the stricter notion of inverse, adjoint functors supply a powerful tool that can capture relations too subtle for the blunt instrument of inverses (which really only “care” about whether the underlying structures are the *same*).

Conceptually, the role adjunctions play in the *theory* of category theory is not dissimilar to what can be seen in the philosophical idea—developed in very different contexts by a variety of different thinkers, each with their own different aims and terminology—of *dialectics*. Proponents of dialectics essentially start out by maintaining that—very abstractly speaking—any notion or model of the “space of concepts” as a sort of static inventory of individual concepts, each set off on its own and capable of “determining itself,” is not to be taken at face value. Unpacking an individual concept by using certain abstract transformations or determinations that lead it to be rearticulated against the backdrop of a concept other than it, then doing the same for the latter concept, pushing it into the setting of the former—these transformations will help to extract previously unarticulated antagonisms in the concepts and uncover ways the original concepts were not as capable of “determining themselves” as they first appeared. Finally, combining the two determinations into a sort of conceptual “unity,” gives a new and more dynamic understanding of the original concepts, each revealed in some fashion as being *what it is* by the way it achieves this unity of distinct determinations in relation to a concept other than it. Using this very schematic idea, one might think of adjunctions as similarly yielding generalized (conceptual) “unities,” where this weakened notion is meant to name the “closest thing” to an *equivalence*.

Philosophers (like Hegel) who are most associated with the development of a concept of dialectics were largely concerned with one particular determination—namely, *negation*. To the extent that we could compare the heuristic supplied by the concept of dialectics and the precise notion of adjunction, we might thus say such philosophers were largely bound by consideration of (proto)adjunctions involving the negation functor and categories in relation to their “opposing” categories. But these situations form a very small, and extreme, slice of the sorts of adjoint relations we can find. In this way, one might say that the notion of adjunction is not just a way of making precise the imprecise idea of dialectics but also a way of releasing it from its arbitrary restriction to involving only functors of the “negation” type.

In short, you could think of the notion of an adjunction as representing something like the wedding of the category-theoretic notion of *universality* with (a much improved version of) the conceptual notion of a *dialectic* or back-and-forth between self-determination and determination-by-another. In chapter 11, we will return to this idea and exhibit adjoint triples that make the connections with the notion of dialectics more precise.

This is a section of [doi:10.7551/mitpress/12581.001.0001](https://doi.org/10.7551/mitpress/12581.001.0001)

Sheaf Theory through Examples

By: Daniel Rosiak

Citation:

Sheaf Theory through Examples

By: Daniel Rosiak

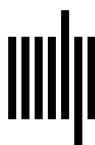
DOI: [10.7551/mitpress/12581.001.0001](https://doi.org/10.7551/mitpress/12581.001.0001)

ISBN (electronic): 9780262370424

Publisher: The MIT Press

Published: 2022

The open access edition of this book was made possible by generous funding and support from Arcadia – a charitable fund of Lisbet Rausing and Peter Baldwin, and MIT Press Direct to Open



The MIT Press

© 2022 Massachusetts Institute of Technology

This work is subject to a Creative Commons CC-BY-NC-ND license.

Subject to such license, all rights are reserved.



The open access edition of this book was made possible by generous funding from Arcadia—a charitable fund of Lisbet Rausing and Peter Baldwin.



The MIT Press would like to thank the anonymous peer reviewers who provided comments on drafts of this book. The generous work of academic experts is essential for establishing the authority and quality of our publications. We acknowledge with gratitude the contributions of these otherwise uncredited readers.

This book was set in LaTeX by the author.

Library of Congress Cataloging-in-Publication Data

Names: Rosiak, Daniel, author.

Title: Sheaf theory through examples / Daniel Rosiak.

Description: Cambridge, Massachusetts : The MIT Press, [2022] | Includes bibliographical references and index.

Identifiers: LCCN 2021058949 | ISBN 9780262542159 (paperback)

Subjects: LCSH: Sheaf theory.

Classification: LCC QA612.36 .R67 2022 | DDC 514/.224—dc23/eng20220521

LC record available at <https://lccn.loc.gov/2021058949>