

CONSTRUCTIVITY IN MATHEMATICS

by

John L. Bell

Contents

1. Intuitionism and Constructive Reasoning

2. A Constructive Look at the Real Numbers

CONSTRUCTIVE MEANING OF THE LOGICAL OPERATORS ORDER ON \mathbb{R}

A CONSTRUCTIVE VERSION OF CANTOR'S THEOREM

ALGEBRAIC OPERATIONS ON \mathbb{R}

CONVERGENCE OF SEQUENCES AND COMPLETENESS OF \mathbb{R}

FUNCTIONS ON \mathbb{R}

3. Intuitionistic Logic

INTUITIONISTIC LOGIC AS A NATURAL DEDUCTION SYSTEM

KRIPKE SEMANTICS AND THE COMPLETENESS THEOREM

THE DISJUNCTION PROPERTY

INTUITIONISTIC LOGIC IN LINEAR STYLE

HEYTING ALGEBRAS AND ALGEBRAIC INTERPRETATIONS OF

INTUITIONISTIC LOGIC

INTUITIONISTIC FIRST-ORDER ARITHMETIC

4. Interlude: Constructivity in Mathematics before

Brouwer

5. Intuitionistic Set Theory

INTUITIONISTIC ZERMELO SET THEORY

DEFINITIONS OF “FINITE”

INTUITIONISTIC ZERMELO-FRAENKEL SET THEORY: ORDINALS

6. Smooth Infinitesimal Analysis

ALGEBRAIC AND ORDER STRUCTURE OF \mathbf{R}

SIA VERSUS CONSTRUCTIVE ANALYSIS

INDECOMPOSABILITY OF THE CONTINUUM IN **SIA**

NATURAL NUMBERS AND INVERTIBLE INFINITESIMALS IN **SIA**

APPLICATIONS OF **SIA** TO PHYSICS

1. Intuitionism and Constructive Reasoning

Intuitionism is the creation of *L. E. J. Brouwer* (1882-1966). Like Kant, Brouwer believed that mathematical concepts are admissible only if they are adequately grounded in *intuition*, that mathematical theories are significant only if they concern entities which are constructed out of something given immediately in intuition, that mathematical definitions must always be constructive, and that the completed infinite is to be rejected. Thus, like Kant, Brouwer held that mathematical theorems are synthetic *a priori* truths. In *Intuitionism and Formalism* (1912), while admitting that the emergence of noneuclidean geometry had discredited Kant's view of space, he maintained, in opposition to the logicians (whom he called "formalists") that arithmetic, and so all mathematics, must derive from the *intuition of time*. In his own words:

Neointuitionism considers the falling apart of moments of life into qualitatively different parts, to be reunited only while remaining separated by time, as the fundamental phenomenon of the human intellect, passing by abstracting from its emotional content into the fundamental phenomenon of mathematical thinking, the intuition of the bare two-oneness. This intuition of two-oneness, the basal intuition of mathematics, creates not only the numbers one and two, but also all finite ordinal numbers, inasmuch as one of the elements of the two-oneness may be thought of as a new two-oneness, which process may be repeated indefinitely; this gives rise still further to the smallest infinite ordinal ω . Finally

this basal intuition of mathematics, in which the connected and the separate, the continuous and the discrete are united, gives rise immediately to the intuition of the linear continuum, i.e., of the “between”, which is not exhaustible by the interposition of new units and which can therefore never be thought of as a mere collection of units. In this way the apriority of time does not only qualify the properties of arithmetic as synthetic a priori judgments, but it does the same for those of geometry, and not only for elementary two- and three-dimensional geometry, but for non-euclidean and n-dimensional geometries as well. For since Descartes we have learned to reduce all these geometries to arithmetic by means of coordinates.

For Brouwer, intuition meant essentially what it did to Kant, namely, the mind’s apprehension of what it has itself constructed; on this view, the only acceptable mathematical proofs are *constructive*. A constructive proof may be thought of as a kind of “thought experiment” —the performance, that is, of an experiment in imagination. According to *Arend Heyting* (1898–1980), a leading member of the intuitionist school,

Intuitionistic mathematics consists ... in mental constructions; a mathematical theorem expresses a purely empirical fact, namely, the success of a certain construction. “ $2 + 2 = 3 + 1$ ” must be read as an abbreviation for the statement “I have effected the mental construction indicated by ‘ $2 + 2$ ’ and ‘ $3 + 1$ ’ and I have found that they lead to the same result.”

From passages such as these one might infer that for intuitionists mathematics is a purely subjective activity, a kind of introspective

reportage, and that each mathematician has a personal mathematics. Certainly they reject the idea that mathematical thought is dependent on any special sort of language, even, occasionally, claiming that, at bottom, mathematics is a “languageless activity”. Nevertheless, the fact that intuitionists evidently regard mathematical theorems as being valid for all intelligent beings indicates that for them mathematics has, if not an objective character, then at least a *transsubjective* one.

A major impact of the intuitionists’ program of constructive proof has been in the realm of *logic*. Brouwer maintained, in fact, that the applicability of traditional logic to mathematics

was caused historically by the fact that, first, classical logic was abstracted from the mathematics of the subsets of a definite finite set, that, secondly, an a priori existence independent of mathematics was ascribed to the logic, and that, finally, on the basis of this supposed apriority it was unjustifiably applied to the mathematics of infinite sets.

Thus Brouwer held that much of modern mathematics is based, not on sound reasoning, but on an illicit extension of procedures valid only in the restricted domain of the finite. He therefore embarked on the heroic course of setting the whole of existing mathematics aside and starting afresh, using only concepts and modes of inference that could be given clear intuitive justification. He hoped that, once enough of the program had been carried out, one could discern the logical laws that intuitive, or constructive, mathematical reasoning actually obeys, and so be able to compare the resulting *intuitionistic*, or *constructive*, *logic*¹ with classical logic.

¹This is not to say that Brouwer was primarily interested in *logic*, far from it: indeed, his distaste for formalization led him not to take very seriously subsequent codifications of intuitionistic logic.

As we have already seen, in constructive mathematical reasoning *an existential statement can be considered affirmed only when an instance is produced*,² and *a disjunction can be considered affirmed only when an explicit one of the disjuncts is demonstrated*. Consequently, neither the classical law of excluded middle³ nor the law of strong reductio ad absurdum⁴ are constructively acceptable. These conclusions have already been noted in connection with the real numbers; let us employ some straightforward examples involving the natural numbers to draw the same conclusions more simply.

Consider the existential statement *there exists an odd perfect number* (i.e., an odd number equal to the sum of its proper divisors) which we shall write as $\exists nP(n)$. Its contradictory is the statement $\forall n\neg P(n)$. Classically, the law of excluded middle then allows us to affirm the disjunction

$$\exists nP(n) \vee \forall n\neg P(n) \tag{1}$$

Constructively, however, in order to affirm this disjunction we must *either* be in a position to affirm the first disjunct $\exists nP(n)$, i.e., to possess, or have the means of obtaining, an odd perfect number, *or* to affirm the second disjunct $\forall n\neg P(n)$, i.e. to possess a demonstration that no odd number is perfect. Since at the present time mathematicians have neither of these⁵, the disjunction (1), and *a fortiori* the law of excluded middle is not constructively admissible.

² Hermann Weyl said of nonconstructive existence proofs that “they inform the world that a treasure exists without disclosing its location.”

³This is the assertion that, for any proposition p , either p or its negation $\neg p$ holds.

⁴This is the assertion that, for any proposition p , $\neg\neg p$ implies p .

⁵And indeed may never have; for little if any progress has been made on the ancient problem of the existence of odd perfect numbers.

It might be thought that, if in fact the second disjunct in (1) is *false*, that is, not every number falsifies P , then we can actually find a number satisfying P by the familiar procedure of testing successively each number $0, 1, 2, 3, \dots$ and breaking off when we find one that does: in other words, that from $\neg\forall n\neg P(n)$ we can infer $\exists nP(n)$. Classically, this is perfectly correct, because the *classical* meaning of $\neg\forall n\neg P(n)$ is “ $P(n)$ will not as a matter of *fact* be found to fail for every number n .” But *constructively* this latter statement has no meaning, because it presupposes that every natural number *has already been constructed* (and checked for whether it satisfies P). Constructively, the statement must be taken to mean something like “we can derive a contradiction from the supposition that we could prove that $P(n)$ failed for every n .” From this, however, we clearly cannot extract a guarantee that, by testing each number in turn, we shall eventually find one that satisfies P . So we see, once again, that the law of strong reductio ad absurdum also fails to be constructively admissible.

As a simple example of a classical existence proof which fails to meet constructive standards, consider the assertion

there exists a pair of irrational real numbers a, b such that a^b is rational.

Classically, this can be proved as follows: let $b = \sqrt{2}$; then b is irrational. If b^b is rational, let $a = b$; then we are through. If b^b is irrational, put $a = b^b$; then $a^b = 2$, which is rational. But in this proof we have not *explicitly identified* a ; we do not know, in fact, whether $a = \sqrt{2}$ or⁶ $a = \sqrt{2}^{\sqrt{2}}$, and it is therefore constructively unacceptable.

⁶In fact a much deeper argument shows that $2^{\sqrt{2}}$ is irrational, and is therefore the correct value of a .

Constructive reasoning differs from its classical counterpart in that it attaches a stronger meaning to some of the logical operators. It has become customary, following Heyting, to explain this stronger meaning in terms of the primitive relation *α is a proof of p* , between mathematical constructions α and mathematical assertions p . To assert the *truth* of p is to assert that one has a construction α such that α is a proof of p ⁷. The meaning of the various logical operators in this scheme is spelt out by specifying how proofs of composite statements depend on proofs of their constituents. Thus:

1. α is a proof of $p \wedge q$ means: α is a pair (β, γ) consisting of a proof β of p and γ of q .
2. α is a proof of $p \vee q$ means: α is a pair (n, β) consisting of a natural number n and a construction β such that, if $n = 0$, then β is a proof of p , and if $n \neq 0$, then β is a proof of q .
3. α is a proof of $p \rightarrow q$ means: α is a construction that converts any proof of p into a proof of q ;
4. α is a proof of $\neg p$ means: α is a construction that shows that no proof of p is possible.

In order to deal with quantified statements we assume that some domain of individuals D is given. Then

5. α is a proof of $\exists x p(x)$ means: α is a pair (d, β) where d is a specified member of D and β is a proof that $p(d)$.

⁷Here by *proof* we are to understand a mathematical construction that establishes the assertion in question, *not* a derivation in some formal system. For example, a proof of $2 + 3 = 5$ in this sense consists of successive constructions of 2, 3 and 5, followed by a construction that adds 2 and 3, finishing up with a construction that compares the result of this addition with 5.

6. α is a proof of $\forall x p(x)$ means: α is a construction which, applied to any member d of D , yields a proof $\alpha(d)$ of $p(d)$.

It is readily seen that, for example, the law of excluded middle is not generally true under this ascription of meaning to the logical operators. For a proof of $p \vee \neg p$ is a pair (β, n) in which c is either a proof of p or a construction showing that no proof of p is possible, and there is nothing inherent in the concept of mathematical construction that guarantees, for an arbitrary proposition p , that either will ever be produced.

As shown by Gödel in the 1930s, it is possible to represent the strengthened meaning of the constructive logical operators in a classical system augmented by the concept of *provability*. If we write $\Box p$ for “ p is provable”, then the scheme below correlates constructive statements with their classical translates.

<i>Constructive</i>	<i>Classical</i>
$\neg p$	$\Box \neg \Box p$
$p \wedge q$	$\Box p \wedge \Box q$
$p \vee q$	$\Box p \vee \Box q$
$p \rightarrow q$	$\Box(\Box p \rightarrow \Box q)$

The translate of the sentence $p \vee \neg p$ is then $\Box p \vee \Box \neg \Box p$, which is (assuming $\Box \Box p \leftrightarrow \Box p$) equivalent to $\neg \Box p \rightarrow \Box \neg \Box p$, that is, to the assertion

if p is not provable, then it is provable that p is not provable.

The fact that there is no *a priori* reason to accept this “solubility” principle lends further support to the intuitionists’ rejection of the law of excluded middle.

Another interpretation of constructive reasoning is provided by *Kolmogorov’s calculus of problems* (A. N. Kolmogorov, 1903–1987). If we denote problems by letters and $a \wedge b$, $a \vee b$, $a \rightarrow b$, $\neg a$ are construed respectively as the problems

to solve both a and b

to solve at least one of a and b

to solve b, given a solution of a

to deduce a contradiction from the hypothesis that a is solved,

then a formal calculus can be set up which coincides with the constructive logic of propositions.

2. A Constructive Look at the Real Numbers.

In constructive mathematics, a problem is counted as solved only if an explicit solution can, in principle at least, be produced. Thus, for example, “There is an x such that $P(x)$ ” means that, in principle at least, we can explicitly produce an x such that $P(x)$. If the solution to the problem involves parameters, we must be able to present the solution explicitly by means of some *algorithm* or *rule* when give values of the parameters. That is, “for every x there is a y such that $P(x, y)$ ” means that, we possess an explicit method of determining, for any given x , a y for which $P(x, y)$. This leads us to examine what it means for a mathematical object to be explicitly given. To begin with, everybody knows what it means to give an *integer* explicitly. For example, $7 \cdot 10^4$ is given explicitly, while the number n defined to be 0 if an odd perfect number exists, and 1 if an odd perfect number does not exist, is not given explicitly. The number of primes less than, say, $10^{1000000}$ is given explicitly, in the sense intended here, since we could, *in principle at least*, calculate this number. Constructive mathematics as we shall understand it is not concerned with questions of feasibility, nor in particular with what can actually be computed in real time by actual computers.

Rational numbers may be defined as pairs of integers (a, b) without a common divisor (where $b > 0$ and a may be positive or negative, or a is 0 and b is 1). The usual arithmetic operations on the rationals, together with the operation of taking the absolute value, are then easily supplied with explicit definitions. Accordingly it is clear what it means to give a rational number explicitly.

To specify exactly what is meant by giving a *real number* explicitly is not quite so simple. For a real number is by its nature an infinite object, but one normally regards only finite objects as capable of being given explicitly. We shall get round this difficulty by stipulating that, to be given a real number, we must be given a (finite) *rule* or *explicit procedure* for calculating it to any desired degree of accuracy. Intuitively speaking, to be given a real number r is to be given a method of computing, for each positive integer n , a rational number r_n such that

$$|r - r_n| < 1/n.$$

These r_n will then obey the law

$$|r_m - r_n| \leq 1/m + 1/n.$$

So, given any numbers k, p , we have, setting $n = 2k$,

$$|r_{n+p} - r_n| \leq 1/(n+p) + 1/n \leq 2/n = 1/k.$$

We are thus led to *define* a real number to be a sequence of rationals $(r_n) = r_1, r_2, \dots$ such that, for any k , a number n can be found such that

$$|r_{n+p} - r_n| \leq 1/k \text{ for all } p.$$

Here we understand that to be given a *sequence* we must be in possession of a *rule* or explicit method for generating its members. Each rational number α may be regarded as a real number by

identifying it with the real number (α, α, \dots) . The set of all real numbers will be denoted, as usual, by \mathbb{R} .

Now of course, for any “given” real number there are a variety of ways of giving explicit approximating sequences for it. Thus it is necessary to define an equivalence relation, “equality on the reals”. The correct definition here is: $r =_R s$ iff for any k , a number n can be found so that

$$|r_{n+p} - s_{n+p}| \leq 1/k \text{ for all } p.$$

When we say that two real numbers are equal we shall mean that they are equivalent in this sense, and so write simply “=” for “ $=_R$ ”

CONSTRUCTIVE MEANING OF THE LOGICAL OPERATORS

It is appropriate here to make a few remarks on the *constructive meaning of the logical operators*. To begin with, if the symbol “ \exists ” is taken to mean “explicit existence” in the sense described above, it cannot be expected to obey the laws of classical logic. For example, $\neg\forall$ is classically equivalent to $\exists\neg$, but the mere knowledge that something cannot always occur does not enable us actually to *determine* a location where it fails to occur. This is generally the case with existence proofs by contradiction. For instance, consider the following standard proof of the *Fundamental Theorem of Algebra*: every polynomial p of degree > 0 has a (complex) zero. If p lacks a zero, then $1/p$ is entire and bounded, and so by Liouville’s theorem must be constant. This proof gives no hint of how actually to construct a zero. (But constructive proofs of this theorem are known.)

The constructive meaning of disjunction is given by the equivalence

$$A \vee B \Leftrightarrow \exists n[(n = 0 \rightarrow A) \& (n \neq 0 \rightarrow B)].$$

That is, $A \vee B$ means that one of A or B holds, and *we can tell which one*.

The constructive meaning of negation is simple: $\neg A$ means that *A leads to a contradiction*. Combining this with the meaning of disjunction enables us to grasp the constructive meaning of the *law of excluded middle*: $A \vee \neg A$ is now seen to express the nontrivial claim that we have a method of deciding which of A or $\neg A$ holds, that is, a method of either proving A or deducing a contradiction from A . If A is an unsolved problem, this claim is dubious at best.

Is it constructively true, for instance, that for any real numbers x and y , we have $x = y \vee x \neq y$? As we shall see, the answer is no. If this assertion were constructively true, then, in particular, we would have a method of deciding whether, for any given rational number r , whether $r = \pi^{\sqrt{2}}$ or not. But at present no such method is known—it is not known, in fact, whether $\pi^{\sqrt{2}}$ is rational or irrational. We can, of course, calculate $\pi^{\sqrt{2}}$ to as many decimal places as we please, and if in actuality it is unequal to a given rational number r , we shall discover this fact after a sufficient amount of calculation. If, however, $\pi^{\sqrt{2}}$ is *equal* to r , even several centuries of computation cannot make this fact certain; we can be sure only that it is very close to r . We have no method which will tell us, in finite time, whether $\pi^{\sqrt{2}}$ exactly coincides with r or not.

This situation may be summarized by saying that equality on the reals is *not decidable*. (By contrast, equality on the integers or rational numbers is decidable.) Observe that this does *not* mean $\neg(x = y \vee x \neq y)$. We have not actually derived a *contradiction* from the assumption $x = y \vee x \neq y$, we have only given an example showing its implausibility. It is natural to ask whether it can actually be *refuted*. For this it would be necessary to make some assumption concerning the real numbers which contradicts classical mathematics. Certain schools of constructive mathematics are willing to make such assumptions; but the majority of constructivists confine themselves to methods which are also classically correct. (Later on, however, we shall describe a model of the real line in which the decidability of equality can be refuted.)

Despite the fact that equality of real numbers is not a decidable relation, it is *stable* in the sense of satisfying the *law of double negation* $\neg(r \neq s) \Rightarrow r = s$. For, given k , we may choose n so that $|r_{n+p} - r_n| \leq 1/4k$ and $|s_{n+p} - s_n| \leq 1/4k$ for all p . If $|r_n - s_n| \geq 1/k$, then we would have $|r_{n+p} - s_{n+p}| \geq 1/2k$ for all p , which entails $r \neq s$. If $\neg(r \neq s)$, it follows that $|r_n - s_n| < 1/k$ and $|r_{n+p} - s_{n+p}| \leq 2/k$ for every p . Since for every k we can find n so that this inequality holds for every p , it follows that $r = s$.

One should not, however, conclude from the stability of equality that the law of double negation $\neg\neg A \rightarrow A$ is generally affirmable. That it is not so can be seen from the following example. Write the decimal expansion of π and below the decimal expansion $\rho = 0.333\dots$, terminating it as soon as a sequence of digits 0123456789 has appeared in π . Then if the 9 of the first sequence 0123456789 in π is the k^{th} digit after the decimal point, $r = (10^k - 1)/3 \cdot 10^k$. Now suppose that ρ were not rational; then $r =$

$(10^k - 1)/3 \cdot 10^k$ would be impossible and no sequence 0123456789 could appear in π , so that $\rho = 1/3$, which is also impossible. Thus the assumption that ρ is not rational leads to a contradiction; yet we are not warranted to assert that ρ is rational, for this would mean that we could calculate integers m and n for which $\rho = m/n$. But this evidently requires that we can produce a sequence 0123456789 in π or demonstrate that no such sequence can appear, and at present we can do neither.

To assert the inequality of two real numbers is constructively weak. In constructive mathematics a stronger notion of inequality, that of *apartness*, is normally used instead. We say that r and s are *apart*, written $r \neq s$, if n and k can actually be found so that $|r_{n+p} - s_{n+p}| > 1/k$ for all p . Clearly $r \neq s$ implies $r \neq s$, but the converse cannot be affirmed constructively.⁸ The proof of $\neg r \neq s \Rightarrow r = s$ given above actually establishes something stronger, namely $\neg r \neq s \Rightarrow r = s$.

ORDER ON \mathbb{R}

The *order relation* on the reals is given constructively by stipulating that $r < s$ is to mean that we have an explicit lower bound on the distance between r and s . That is,

$$r < s \Leftrightarrow n \text{ and } k \text{ can be found so that } s_{n+p} - r_{n+p} > 1/k \text{ for all } p.$$

It can readily be shown that, for any real numbers x, y such that $x < y$, there is a rational number α such that $x < \alpha < y$.

⁸ In fact the converse is equivalent to *Markov's Principle*, which asserts that, if, for each n , $x_n = 0$ or 1 , and if it is contradictory that $x_n = 0$ for all n , then there exists n for which $x_n = 1$. This thesis is accepted by some, but not all schools of constructivism.

We observe that $r \neq s \Leftrightarrow r < s \vee s < r$. The implication from right to left is clear. Conversely, suppose that $r \neq s$. Find n and k so that $|r_{n+p} - s_{n+p}| > 1/k$ for every p , and determine $m > n$ so that $|r_m - r_{m+p}| < 1/4k$ and $|s_m - s_{m+p}| < 1/4k$ for every p . Either $r_m - s_m > 1/k$ or $s_m - r_m > 1/k$; in the first case $r_{m+p} - s_{m+p} > 1/2k$ for every p , whence $s < r$; similarly, in the second case, we obtain $r < s$.

We define $r \leq s$ to mean that $s < r$ is false. Notice that $r \leq s$ is not the same as $r < s$ or $r = s$: in the case of the real number ρ defined above, for instance, clearly $\rho \leq 1/3$, but we do not know whether $\rho < 1/3$ or $\rho = 1/3$. Still, it is true that $r \leq s \wedge s \leq r \Rightarrow r = s$. For the premise is the negation of $r < s \vee s < r$, which, by the above, is equivalent to $\neg r \neq s$. But we have already seen that this last implies $r = s$.

There are several common properties of the order relation on real numbers which hold classically but which cannot be established constructively. Consider, for example, the trichotomy law $x < y \vee x = y \vee y < x$. Suppose we had a method enabling us to decide which of the three alternatives holds. Applying it to the case $y = 0$, $x = \pi^{\sqrt{2}} - r$ for rational r would yield an algorithm for determining whether $\pi^{\sqrt{2}} = r$ or not, which we have already observed is an open problem. One can also demonstrate the failure of the trichotomy law (as well as other classical laws) by the use of “fugitive sequences”. Here one picks an unsolved problem of the form $\forall n P(n)$, where P is a decidable property of integers—for example, Goldbach’s conjecture that every even number ≥ 4 is the sum of two odd primes. Now one defines a sequence—a “fugitive” sequence—of integers (n_k) by $n_k = 0$ if $2k$ is the sum of two primes and $n_k = 1$ otherwise. Let r be the real number defined by $r_k = 0$ if $n_k = 0$ for all $j \leq k$, and $r_k = 1/m$ otherwise, where m is the least positive integer such that $n_m = 1$. It is then easy to check that $r \geq 0$

and $r = 0$ iff Goldbach's conjecture holds. Accordingly the correctness of the trichotomy law would imply that we could resolve Goldbach's conjecture. Of course, Goldbach's conjecture might be resolved in the future, in which case we would merely choose another unsolved problem of a similar form to define our fugitive sequence.

A similar argument shows that the law $r \leq s \vee s \leq r$ also fails constructively: define the real number s by $s_k = 0$ if $n_k = 0$ for all $j \leq k$; $s_k = 1/m$ if m is the least positive integer such that $n_m = 1$, and m is even; $s_k = -1/m$ if m is the least positive integer such that $n_m = 1$, and m is odd. Then $s \leq 0$ (resp. $0 \leq s$) would mean that there is no number of the form $2 \cdot 2k$ (resp. $2 \cdot (2k + 1)$) which is not the sum of two primes. Since neither claim is at present known to be correct, we cannot assert the disjunction $s \leq 0 \vee 0 \leq s$.

In constructive mathematics there is a convenient substitute for trichotomy known as the *comparison principle*. This is the assertion

$$r < t \Rightarrow r < s \vee s < t.$$

Its validity can be established in a manner similar to the foregoing.

A CONSTRUCTIVE VERSION OF CANTOR'S THEOREM

Cantor's theorem that \mathbb{R} is uncountable has the following constructive version:

Theorem. Let (a_n) be a sequence of real numbers, and let x_0 and y_0 be real numbers with $x_0 < y_0$. Then there exists a real number x such that $x_0 \leq x \leq y_0$ and $x \neq a_n$ for all $n \geq 1$.

Proof. We construct by recursion sequences (x_n) , (y_n) of rational numbers such that

- (i) $x_0 \leq x_n \leq x_m < y_m \leq y_n \leq y_0$ ($m \geq n \geq 1$)
- (ii) $x_n > a_n$ or $y_n < a_n$ ($n \geq 1$)
- (iii) $y_n - x_n < n^{-1}$ ($n \geq 1$).

Assume that $n \geq 1$ and that $x_0, \dots, x_{n-1}, y_0, \dots, y_{n-1}$ have been constructed. Either $a_n > x_{n-1}$ or $a_n < y_{n-1}$. If the former, let x_n be any rational number with $x_{n-1} < x_n < \min(a_n, y_{n-1})$ and let y_n be any rational number with $x_n < y_n < \min(a_n, y_{n-1}, x_n + 1/n)$. The relevant inequalities are then satisfied. If $a_n < y_{n-1}$, let y_n be any rational number with $\max(a_n, x_{n-1}) < y_n < y_{n-1}$ and let x_n be any rational number with $\max(a_n, x_{n-1}, y_n - 1/n) < x_n < y_n$. The relevant inequalities are again satisfied.

From (i) and (iii) it follows that

$$|x_m - x_n| = x_m - x_n < y_m - x_n < 1/n \quad (m \geq n)$$

Similarly $|y_m - y_n| < 1/n$ for $m \geq n$. Therefore $x = (x_n)$ and $y = (y_n)$ are real numbers. By (i) and (iii), they are equal. By (i), $x_n \leq x$ and $y_n \geq y$ for all n . If $a_n < x_n$, then $a_n < x$ and so $a_n \neq x$; if $a_n > y_n$, then $a_n > y = x$ and again $a_n \neq x$. Accordingly x has the required properties. #

The fundamental operations $+$, $-$, \cdot , $^{-1}$ and $| \cdot |$ are defined for real numbers as one would expect, viz.

- $r + s$ is the sequence $(r_n + s_n)$
- $r - s$ is the sequence $(r_n - s_n)$
- $r \cdot s$ or rs is the sequence $(r_n s_n)$
- if $r \neq 0$, r^{-1} is the sequence (t_n) , where $t_n = r_n^{-1}$ if $t_n \neq 0$ and $t_n = 0$ if $r_n = 0$
- $|r|$ is the sequence $(|r_n|)$

It is then easily shown that $rs \neq 0 \Leftrightarrow r \neq 0 \wedge s \neq 0$. For if $r \neq 0 \wedge s \neq 0$, we can find k and n such that $|r_{n+p}| > 1/k$ and $|s_{n+p}| > 1/k$ for every p , so that $|r_{n+p}s_{n+p}| > 1/k^2$ for every p , and $rs \neq 0$. Conversely, if $rs \neq 0$, then we can find k and n so that

$$|r_{n+p}s_{n+p}| > 1/k, \quad |r_{n+p} - r_n| < 1, \quad |s_{n+p} - s_n| < 1$$

for every p . It follows that

$$|r_{n+p}| > 1/k(|s_n| + 1) \text{ and } |s_{n+p}| > 1/k(|r_n| + 1)$$

for every p , whence $r \neq 0 \wedge s \neq 0$.

But it is not constructively true that, if $rs = 0$, then $r = 0$ or $s = 0$! To see this, use the following prescription to define two real numbers r and s . If in the first n decimals of π no sequence 0123456789 occurs, put $r_n = s_n = 2^{-n}$; if a sequence of this kind does occur in the first n decimals, suppose the 9 in the first such sequence is the k^{th} digit. If k is odd, put $r_n = 2^{-k}$, $s_n = 2^{-n}$; if k is even, put $r_n = 2^{-n}$, $s_n = 2^{-k}$. Then we are unable to decide whether $r = 0$ or $s = 0$. But $rs = 0$. For in the first case above $r_n s_n = 2^{-2n}$; in the

second $r_n s_n = 2^{-k-n}$. In either case $|r_n s_n| < 1/m$ for $n > m$, so that $rs = 0$.

CONVERGENCE OF SEQUENCES AND COMPLETENESS OF \mathbb{R}

As usual, a sequence (a_n) of real numbers is said to *converge* to a real number b , or to have *limit* s if, given any natural number k , a natural number n can be found so that for every natural number p ,

$$|b - a_{n+p}| < 2^{-k}.$$

As in classical analysis, a constructive necessary and condition that a sequence (a_n) of real numbers be convergent is that it be a *Cauchy* sequence, that is, if, given any given any natural number k , a natural number n can be found so that for every natural number p ,

$$|a_{n+p} - a_n| < 2^{-k}.$$

But some classical theorems concerning convergent sequences are no longer valid constructively. For example, a bounded momotone sequence need no longer be convergent. A simple counterexample is provided by the sequence (a_n) defined as follows: $a_n = 1 - 2^{-n}$ if among the first n digits in the decimal expansion of π no sequence 0123456789 occurs, while $a_n = 2 - 2^{-n}$ if among these n digits such a sequence does occur. Since it is not known whether the limit of this sequence, if it exists, is 1 or 2, we cannot claim that that this limit exists as a well defined real number.

In classical analysis \mathbb{R} is *complete* in the sense that every nonempty set of real numbers that is bounded above has a supremum. As it stands, this assertion is constructively incorrect. For consider the set A of members $\{x_1, x_2, \dots\}$ of any fugitive sequence of 0s and 1s. Clearly A is bounded above, and its supremum would be either 0 or 1. If we knew which, we would also know whether $x_n = 0$ for all n , and the sequence would no longer be fugitive.

However, the completeness of \mathbb{R} can be salvaged by defining suprema and infima somewhat more delicately than is customary in classical mathematics. A nonempty set A of real numbers is *bounded above* if there exists a real number b , called an *upper bound* for A , such that $x \leq b$ for all $x \in A$. A real number b is called a *supremum*, or *least upper bound*, of A if it is an upper bound for A and if for each $\varepsilon > 0$ there exists $x \in A$ with $x > b - \varepsilon$. We say that A is *bounded below* if there exists a real number b , called a *lower bound* for A , such that $b \leq x$ for all $x \in A$. A real number b is called an *infimum*, or *greatest lower bound*, of A if it is a lower bound for A and if for each $\varepsilon > 0$ there exists $x \in A$ with $x < b + \varepsilon$. The supremum (respectively, infimum) of A , is unique if it exists and is written $\sup A$ (respectively, $\inf A$).

We now prove the *constructive least upper bound principle*.

Theorem. Let A be a nonempty set of real numbers that is bounded above. Then $\sup A$ exists if and only if for all $x, y \in \mathbb{R}$ with $x < y$, either y is an upper bound for A or there exists $a \in A$ with $x < a$.

Proof. If $\sup A$ exists and $x < y$, then either $\sup A < y$ or $x < \sup A$; in the latter case we can find $a \in A$ with $\sup A - (\sup A - x) < a$, and hence $x < a$. Thus the stated condition is necessary.

Conversely, suppose the stated condition holds. Let a_1 be an element of A , and choose an upper bound b_1 for A with $b_1 > a_1$. We construct recursively a sequence (a_n) in A and (b_n) of upper bounds for A such that, for each $n \geq 0$,

$$(i) \quad a_n \leq a_{n+1} \leq b_{n+1} \leq b_n$$

and

$$(ii) \quad b_{n+1} - a_{n+1} \leq (b_n - a_n).$$

Having found a_1, \dots, a_n and b_1, \dots, b_n , if $a_n + (b_n - a_n)$ is an upper bound for A , put $b_{n+1} = a_n + (b_n - a_n)$ and $a_{n+1} = a_n$; while if there exists $a \in A$ with $a > a_n + (b_n - a_n)$, we set $a_{n+1} = a$ and $b_{n+1} = b_n$. This completes the recursive construction.

From (i) and (ii) we have

$$0 \leq b_n - a_n \leq (b_1 - a_1)^{n-1}.$$

It follows that the sequences (a_n) and (b_n) converge to a common limit ℓ with $a_n \leq \ell \leq b_n$ for $n \geq 1$. Since each b_n is an upper bound for A , so is ℓ . On the other hand, given $\varepsilon > 0$, we can choose n so that $\ell \geq a_n > \ell - \varepsilon$, where $a_n \in A$. Hence $\ell = \sup A$. ■

An analogous result for infima can be stated and proved in a similar way.

FUNCTIONS ON \mathbb{R}

Considered constructively, a *function* from \mathbb{R} to \mathbb{R} is a rule F which enables us, when given a real number x , to compute another real number $F(x)$ in such a way that, if $x = y$, then $F(x) = F(y)$. It is easy to check that every polynomial is a function in this sense, and that various power series and integrals, for example those defining $\tan x$ and e^x , also determine functions.

Viewed constructively, some classically defined “functions” on \mathbb{R} can no longer be considered to be defined on the whole of \mathbb{R} . Consider, for example, the “blip” function B defined by $B(x) = 0$ if $x \neq 0$ and $B(0) = 1$. Here the domain of the function is $\{x \in \mathbb{R} : x = 0 \vee x \neq 0\}$. But we have seen that we cannot assert $\text{dom}(B) = \mathbb{R}$. So the blip function is not well defined as a function from \mathbb{R} to \mathbb{R} . Of course, classically, B is the simplest *discontinuous* function defined on \mathbb{R} . The fact that the simplest possible discontinuous function fails to be defined on the whole of \mathbb{R} gives grounds for the suspicion that *no* function defined on \mathbb{R} can be discontinuous; in other words, that, constructively speaking, *all functions defined on \mathbb{R} are continuous*. (This claim was a central tenet of intuitionism’s founder, Brouwer.) This claim is plausible. For if a function F is well-defined on all reals x , it must be possible to compute the value for all rules x determining real numbers, that is, determining their sequences of rational approximations x_1, x_2, \dots . Now $F(x)$ must be computed to accuracy ε in a finite number of steps—the number of steps depending on ε . This means that only finitely many

approximations can be used, i.e., $F(x)$ can be computed to within ε only when x is known within δ for some δ . Thus F should indeed be continuous. In fact all known examples of constructive functions are continuous.

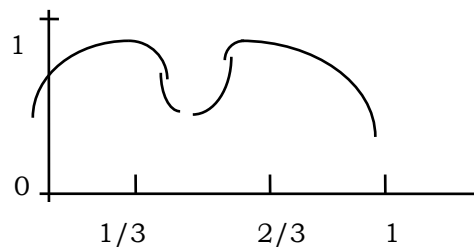
Constructively, a real valued function f is *continuous* if for each $\varepsilon > 0$ there exists $\omega(\varepsilon) > 0$ such that $|f(x) - f(y)| \leq \varepsilon$ whenever $|x - y| < \omega(\varepsilon)$. The operation $\varepsilon \mapsto \omega(\varepsilon)$ is called a *modulus of continuity* for f .

If all functions on \mathbb{R} are continuous, then a subset A of \mathbb{R} may fail to be genuinely *complemented*: that is, there may be no subset B of \mathbb{R} disjoint from A such that $\mathbb{R} = A \cup B$. In fact suppose that A, B are disjoint subsets of \mathbb{R} and that there is a point $a \in A$ which can be approached arbitrarily closely by points of B (or vice-versa). Then, assuming all functions on \mathbb{R} are continuous, it cannot be the case that $\mathbb{R} = A \cup B$. For if so, we may define the function f on \mathbb{R} by $f(x) = 0$ if $x \in A$, $f(x) = 1$ if $x \in B$. Then for all $\delta > 0$ there is $b \in B$ for which $|b - a| < \delta$, but $|f(b) - f(a)| = 1$. So f fails to be continuous at a , and we conclude that $\mathbb{R} \neq A \cup B$.

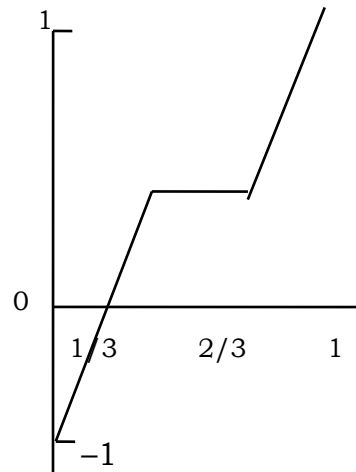
In particular, if we take A to be any finite set of real numbers, any union of open or closed intervals, or the set \mathbb{Q} of rational numbers, then in each case the set B of points “outside” A satisfies the above condition. Accordingly, for each such subset A , \mathbb{R} is not “decomposable” into A and the set of points “outside” A , in the sense that these two sets of points together exhaust \mathbb{R} . This fact indicates that the constructive continuum is a great deal more “cohesive” than its classical counterpart. For classically, the continuum is merely *connected* in the sense that it is not

(nontrivially) decomposable into two open (or closed) subsets. Constructively, however, \mathbb{R} is indecomposable into subsets which are neither open nor closed. Indeed, in some formulations of constructive analysis, \mathbb{R} is cohesive in the ultimate sense that it cannot be decomposed in *any way whatsoever*. In this sense the constructive real line approximates closely to the ideal of a true continuum.

Certain well-known theorems of classical analysis concerning continuous functions fail in constructive analysis. One such is the *theorem of the maximum*: a uniformly continuous function on a closed interval assumes its maximum at some point. For consider, as in the figure below, a function $f : [0,1] \rightarrow \mathbb{R}$ with two relative maxima, one at $x = 1/3$ and the other at $x = 2/3$ and of approximately the same value. Now arrange things so that $f(1/3) = 1$ and $f(2/3) = 1 + t$, where t is some small parameter. If we could tell where f assumes its absolute maximum, clearly we could also determine whether $t \leq 0$ or $t \geq 0$, which, as we have seen, is not, in general, possible. Nevertheless, it can be shown that from f we can in fact calculate the maximum value itself, so that at least one can assert the existence of that maximum, even if one can't tell exactly where it is assume:



Another classical result that fails to hold constructively in its usual form is the well-known *intermediate value theorem*. This is the assertion that, for any continuous function f from the unit interval $[0, 1]$ to \mathbb{R} , such that $f(0) = -1$ and $f(1) = 1$, there exists a real number $a \in [0,1]$ for which $f(a) = 0$. To see that this fails constructively, consider the function f depicted below: here f is piecewise linear, taking the value t (a small parameter) between $x = 1/3$ and $x = 2/3$. If the intermediate value theorem held, we



could determine a for which $f(a) = 0$. Then either $a < 2/3$ or $a > 1/3$; in the former case $t \geq 0$; in the latter $t \leq 0$. Thus we would be able to decide whether $t \geq 0$ or $t \leq 0$; but we have seen that this is not constructively possible in general.

However, it can be shown that, constructively, the intermediate value theorem is “almost” true in the sense that

$$\forall f \forall \varepsilon > 0 \exists a (|f(a)| < \varepsilon)$$

and also in the sense that, if we write $P(f)$ for

$$\forall b \forall a < b \exists c (a < c < b \wedge f(c) \neq 0),$$

then

$$\forall f [P(f) \rightarrow \exists x (f(x) = 0)].$$

This example illustrates how a single classical theorem “refracts” into several constructive theorems.

3. Intuitionistic Logic

INTUITIONISTIC LOGIC AS A NATURAL DEDUCTION SYSTEM

Intuitionistic logic may be elegantly formulated as a *natural deduction system* in a first-order language \mathcal{L} . It will be convenient to omit the negation symbol \neg from \mathcal{L} and introduce instead the falsehood symbol \perp ⁹; $\neg\alpha$ can then be defined as $\alpha \rightarrow \perp$. (We use lower-case Greek letters to denote formulas of \mathcal{L} .) The system here has no axioms, just rules, which are used to generate *derivations*. The simplest rules have the form

$$\frac{\dots\dots\dots}{\alpha}$$

This is to be read: α is an immediate consequence of the premises above the line. Certain rules involve *assumptions* which are later *cancelled*: a cancelled assumption is indicated by putting a cross next to it as in $\times\alpha$.

The rules are of two sorts, introduction rules and elimination rules.

⁹ We conceive of \perp as a “self-contradictory” atomic sentence that has *no* proof. More precisely, \perp is taken to be an “idealised” proposition with the property that each of its proofs can be converted into a proof of *any* proposition whatever.

Introduction rules

$$\boxed{\wedge\mathbf{I} \quad \frac{\alpha \quad \beta}{\alpha \wedge \beta}}$$

$$\boxed{\vee\mathbf{I} \quad \frac{\alpha \quad \beta}{\alpha \vee \beta} \quad \frac{\alpha \quad \beta}{\alpha \vee \beta}}$$

:

:

:

$$\boxed{\rightarrow\mathbf{I} \quad \frac{\begin{array}{|l} \times\alpha \\ \gamma \\ \gamma \end{array}}{\beta}}{\alpha \rightarrow \beta}$$

$$\boxed{\phantom{\rightarrow\mathbf{I}} \quad \frac{}{}}{\alpha}$$

$$\boxed{\forall\mathbf{I} \quad \frac{\alpha(x)}{\forall x \alpha(x)}}$$

$$\boxed{\exists\mathbf{I} \quad \frac{\alpha(y)}{\exists x \alpha(x)}}$$

:

:

β

Elimination rules

$$\boxed{\wedge\mathbf{E} \quad \frac{\alpha \wedge \beta}{\alpha} \quad \frac{\alpha \wedge \beta}{\beta}}$$

$$\boxed{\vee\mathbf{E} \quad \frac{\begin{array}{|l} \times\alpha \\ \\ \end{array}}{\alpha \vee \beta} \quad \gamma}$$

$$\boxed{\phantom{\vee\mathbf{E}} \quad \frac{}{}}{}$$

$$\boxed{\rightarrow\mathbf{E} \quad \frac{\alpha \quad \alpha \rightarrow \beta}{\beta}}$$

$$\boxed{\forall\mathbf{E} \quad \frac{\forall x \alpha(x)}{\alpha(t)}}$$

$$\boxed{\exists\mathbf{E} \quad \frac{\begin{array}{|l} \alpha(t) \\ \alpha(x) \end{array}}{\exists x \alpha(x)} \quad \beta}$$

The quantifier rules are subject to the following conditions: in the rules $\exists\mathbf{I}$ and $\forall\mathbf{E}$, t must be free for x in α . An application of $\forall\mathbf{I}$ is permitted only if the variable x does not occur in any of the assumptions arising in the derivation of $\alpha(x)$, and similarly, in an application of $\exists\mathbf{E}$ the free variable y in the cancelled formula $\alpha(y)$ must not occur free in β or in any of the assumptions in the right-hand derivation of β .

Each of these rules admits easy justification in terms of the constructive meaning of the logical operators spelled out in the previous chapter.

A formula α appearing at the bottom of a derivation D is said to be *derivable* from the (finite) set of uncanceled assumptions in D . If Γ is a set of formulas, we write $\Gamma \vdash \alpha$ to indicate that α is derivable from a subset of Γ . We write $\vdash \alpha$ for $\emptyset \vdash \alpha$ and say that α is *provable*. Here are a couple of derivations to illustrate how provability is established:

$$\alpha \rightarrow \neg\neg\alpha \qquad \frac{\frac{\frac{x^{(1)}\neg\alpha \quad x^{(2)}\alpha}{\perp} \quad \rightarrow\mathbf{E}}{\neg\neg\alpha} \quad \rightarrow\mathbf{I}}{\alpha \rightarrow \neg\neg\alpha} \quad (2)$$

(recall here that $\neg\alpha$ is $\alpha \rightarrow \perp$)

$$\neg\neg\forall x\alpha(x) \rightarrow \forall x\neg\neg\alpha(x) \qquad \frac{\frac{\frac{x^{(1)}\forall x\alpha(x)}{\alpha(x)} \quad x^{(2)}\neg\alpha(x)}{\perp}}{\neg\neg\forall x\alpha(x)} \quad (1) \qquad x^{(3)}\neg\neg\forall x\alpha(x)$$

$$(2) \quad \frac{\frac{\perp}{\neg\neg\alpha(x)}}{\forall x\neg\neg\alpha(x)}$$

$$(3) \quad \neg\neg\forall x\alpha(x) \rightarrow \forall x\neg\neg\alpha(x)$$

Accordingly, $\vdash \alpha \rightarrow \neg\neg\alpha$ and $\vdash \neg\neg\forall x\alpha(x) \rightarrow \forall x\neg\neg\alpha(x)$. Similarly, one can establish the provability of the following formulas:

1. $(\alpha \rightarrow \beta) \rightarrow ((\beta \rightarrow \gamma) \rightarrow (\alpha \rightarrow \gamma))$
2. $(\alpha \rightarrow \beta) \rightarrow (\neg\beta \rightarrow \neg\alpha)$
3. $\neg\alpha \leftrightarrow \neg\neg\neg\alpha$
4. $\neg(\alpha \vee \beta) \leftrightarrow (\neg\alpha \wedge \neg\beta)$
5. $\neg\neg(\alpha \vee \neg\alpha)$
6. $(\alpha \rightarrow \beta) \rightarrow \neg(\alpha \wedge \neg\beta)$
7. $(\alpha \rightarrow \neg\beta) \leftrightarrow \neg(\alpha \wedge \beta)$
8. $(\neg\neg\alpha \wedge \neg\neg\beta) \leftrightarrow \neg\neg(\alpha \wedge \beta)$
9. $(\neg\neg\alpha \rightarrow \neg\neg\beta) \leftrightarrow \neg\neg(\alpha \rightarrow \beta)$
10. $\exists x \neg\alpha(x) \rightarrow \neg\forall x \alpha(x)$
11. $\neg\exists x \alpha(x) \leftrightarrow \forall x \neg \alpha(x)$
12. $\alpha \vee \forall x \beta(x) \rightarrow \forall x (\alpha \vee \beta(x))$
13. $\forall x (\alpha \rightarrow \beta(x)) \leftrightarrow (\alpha \rightarrow \forall x \beta(x))$
14. $\forall x (\alpha(x) \rightarrow \beta) \leftrightarrow (\exists x \alpha(x) \rightarrow \beta)$
15. $\exists x(\alpha \rightarrow \beta(x)) \rightarrow (\alpha \rightarrow \exists x\beta(x))$
16. $\exists x(\alpha(x) \rightarrow \beta) \rightarrow (\forall x \alpha(x) \rightarrow \beta)^*$

Classical logic may be obtained by adding to the rules of intuitionistic logic the rule of (strong) *reductio ad absurdum*, viz.,

$$\begin{array}{l}
 \mathbf{RAA} \qquad \qquad \qquad x \neg \alpha \\
 \qquad \qquad \qquad \qquad \qquad : \\
 \qquad \qquad \qquad \qquad \qquad : \\
 \qquad \qquad \qquad \qquad \qquad \perp \\
 \hline
 \qquad \qquad \qquad \qquad \qquad \alpha
 \end{array}$$

This means that intuitionistic logic is a subsystem of the corresponding classical systems. Nevertheless, as Gödel and Gentzen showed in the 1930s, classical logic can actually be embedded into intuitionistic logic by means of a suitable reinterpretation of classical conjunction and existence. Gödel achieved this by means of his *translation*, assigning to each formula α of \mathcal{L} a formula α^* of \mathcal{L} as follows:

1. $\perp^* = \perp$ and $\alpha^* = \neg\neg\alpha$ for atomic α distinct from \perp
2. $(\alpha \wedge \beta)^* = \alpha^* \wedge \beta^*$
3. $(\alpha \vee \beta)^* = \neg(\neg\alpha^* \wedge \neg\beta^*)$
4. $(\alpha \rightarrow \beta)^* = \alpha^* \rightarrow \beta^*$
5. $(\forall x \alpha(x))^* = \forall x \alpha^*(x)$
6. $(\exists x \alpha(x))^* = \neg\forall x \neg\alpha^*(x)$

Writing Γ^* for $\{\alpha^*: \alpha \in \Gamma\}$, and \vdash_c, \vdash_i for classical and intuitionistic derivability, one proves by induction on derivations that $\Gamma \vdash_c \alpha \Leftrightarrow \Gamma^* \vdash_i \alpha^*$. It follows easily from this that classical predicate (propositional) logic is conservative over intuitionistic predicate (propositional) logic with respect to *negative* formulas, that is, formulas in which all atomic sentence (apart from \perp) occur negated and which contain only the operators $\wedge, \rightarrow, \perp, \forall$. (Observe that

such formulas α satisfy $\vdash_i \alpha^* \leftrightarrow \alpha$.) And we also obtain, for propositional logic, *Glivenko's theorem*: $\vdash_c \alpha \leftrightarrow \vdash_i \neg\neg\alpha^*$. (Observe that, for a propositional formula α , $\vdash_i \alpha \leftrightarrow \alpha^*$.)

KRIPKE SEMANTICS AND THE COMPLETENESS THEOREM

Kripke semantics provides a flexible and suggestive framework for interpreting intuitionistic first-order logic. A *frame* or *Kripke structure* for \mathcal{L} is a quadruple $K = (P, \leq, S)$ where P is a set partially ordered by \leq and S is a function assigning to each element $a \in P$ an \mathcal{L} -structure S_a in such a way that $S_a \subseteq S_b$ whenever $a \leq b$.¹⁰ We say that K is *built on* P . The members of P may be thought of as “stages of knowledge”. We define the relation \Vdash_K of *forcing over* K between members of P and sentences of \mathcal{L} recursively as follows:

- * for atomic σ , $a \Vdash_K \sigma$ if $S_a \models \sigma$ ¹¹
- * $a \Vdash_K \perp$ never
- * $a \Vdash_K \alpha \wedge \beta$ if $a \Vdash_K \alpha$ and $a \Vdash_K \beta$
- * $a \Vdash_K \alpha \vee \beta$ if $a \Vdash_K \alpha$ or $a \Vdash_K \beta$
- * $a \Vdash_K \alpha \rightarrow \beta$ if $\forall b \geq a$ $b \Vdash_K \alpha$ implies $b \Vdash_K \beta$
- * $a \Vdash_K \forall x \alpha(x)$ if $\forall b \geq a \forall u \in |S_b|$ ¹² $b \Vdash_K \alpha(u)$
- * $a \Vdash_K \exists x \alpha(x)$ if $\exists u \in |S_a|$ $a \Vdash_K \alpha(u)$.

¹⁰ If \mathcal{L} is a propositional language, we take S to be a function assigning to each $a \in P$ a set of proposition letters in such a way that $S(a) \subseteq S(b)$ whenever $a \leq b$.

¹¹ When \mathcal{L} is a propositional language this clause becomes: for atomic σ , $a \Vdash_K \sigma$ if $\sigma \in S_a$

¹² Here $|S|$ denotes the domain of a structure S .

Clearly we have

$$* \quad a \Vdash_K \neg\alpha \text{ if } \forall b \geq a \quad b \not\Vdash_K \alpha.$$

Also it is easily shown that

$$a \Vdash_K \neg\neg\alpha \text{ if } \forall b \geq a \exists c \geq b \quad c \Vdash_K \alpha.$$

And by induction one proves that the forcing relation is *persistent*, that is,

$$a \Vdash_K \alpha \ \& \ b \geq a \text{ implies } b \Vdash_K \alpha.$$

Now let Γ be a set of sentences of \mathcal{L} , and K a frame. We write

$$\Vdash_K \alpha \text{ for } \forall a \in P \quad a \Vdash_K \alpha \text{ (here } \alpha \text{ is said to be } \textit{true} \text{ in } K)$$

$$a \Vdash_K \Gamma \text{ for } \forall \alpha \in \Gamma \quad a \Vdash_K \alpha$$

$$\Gamma \Vdash \alpha \text{ for } \forall K \forall a \in P [a \Vdash_K \Gamma \Rightarrow a \Vdash_K \alpha]$$

$$\Vdash \alpha \text{ for } \forall K \quad \Vdash_K \alpha$$

One can now prove the

Soundness Theorem. $\Gamma \vdash \alpha \Rightarrow \Gamma \Vdash \alpha.$

Proof. For simplicity we confine our sketch of a proof of this theorem to the propositional case only. The proof proceeds by induction on the derivation D of α from Γ . We consider the induction steps for the rules $\forall\mathbf{E}$ and $\rightarrow\mathbf{I}$.

$$\begin{array}{c}
 \mathbf{\vee E} \\
 \begin{array}{ccc}
 & \times\alpha & \times\beta \\
 & : & : \\
 & : & : \\
 \alpha \vee \beta & \gamma & \gamma \\
 \hline
 & \gamma &
 \end{array}
 \end{array}$$

Here the induction hypothesis is the conjunction of the following clauses:

$$\forall a [a \Vdash_{\mathcal{K}} \Gamma \Rightarrow a \Vdash_{\mathcal{K}} \alpha \vee \beta], \quad \forall a [a \Vdash_{\mathcal{K}} \Gamma \cup \{\alpha\} \Rightarrow a \Vdash_{\mathcal{K}} \gamma], \quad \forall a [a \Vdash_{\mathcal{K}} \Gamma \cup \{\beta\} \Rightarrow a \Vdash_{\mathcal{K}} \gamma]$$

If $a \Vdash_{\mathcal{K}} \Gamma$ then $a \Vdash_{\mathcal{K}} \alpha$ or $a \Vdash_{\mathcal{K}} \beta$; suppose $a \Vdash_{\mathcal{K}} \alpha$. Then $a \Vdash_{\mathcal{K}} \Gamma \cup \{\alpha\}$ so $a \Vdash_{\mathcal{K}} \gamma$. Similarly when $a \Vdash_{\mathcal{K}} \beta$. Hence $\forall a [a \Vdash_{\mathcal{K}} \Gamma \Rightarrow a \Vdash_{\mathcal{K}} \gamma]$ as required.

$$\begin{array}{c}
 \mathbf{\rightarrow I} \\
 \begin{array}{c}
 \boxed{\alpha} \\
 : \\
 : \\
 \hline \beta \\
 \alpha \rightarrow \beta
 \end{array}
 \end{array}$$

In this case the inductive hypothesis is $\forall a [a \Vdash_{\mathcal{K}} \Gamma \cup \{\alpha\} \Rightarrow a \Vdash_{\mathcal{K}} \beta]$.

We have to establish $\forall a [a \Vdash_{\mathcal{K}} \Gamma \Rightarrow a \Vdash_{\mathcal{K}} \alpha \rightarrow \beta]$, i.e.

$$\forall a [a \Vdash_{\mathcal{K}} \Gamma \Rightarrow \forall b \geq a [b \Vdash_{\mathcal{K}} \alpha \Rightarrow b \Vdash_{\mathcal{K}} \beta]].$$

Suppose that $a \Vdash_{\mathcal{K}} \Gamma$, $b \geq a$, $b \Vdash_{\mathcal{K}} \alpha$. Then $a \Vdash_{\mathcal{K}} \Gamma$ by persistence, so that $b \Vdash_{\mathcal{K}} \Gamma \cup \{\alpha\}$, whence $b \Vdash_{\mathcal{K}} \beta$ by inductive hypothesis, as required. ■

We now set about proving the converse to the soundness theorem, the *completeness theorem*. Again, for simplicity we confine attention to propositional logic.

A *theory* in \mathcal{L} is a set of sentences closed under deducibility. A theory Γ is said to be *prime* if $\perp \notin \Gamma$ and, for any sentences α, β , $\alpha \vee \beta \in \Gamma \Leftrightarrow \alpha \in \Gamma$ or $\beta \in \Gamma$.

Extension Lemma. Suppose $\Gamma \not\vdash \gamma$. Then there is a prime theory Π such that $\Gamma \subseteq \Pi$ and $\gamma \notin \Pi$.

Proof. Enumerate the sentences of \mathcal{L} as $\sigma_0, \sigma_1, \dots$. Define a sequence of sets of sentences $\Gamma_0, \Gamma_1, \dots$ as follows. First, put $\Gamma_0 = \Gamma$. At stage $k + 1$ we distinguish 3 cases.

1. If $\Gamma_k \cup \{\sigma_k\} \vdash \gamma$, put $\Gamma_{k+1} = \Gamma_k$.
2. If $\Gamma_k \cup \{\sigma_k\} \not\vdash \gamma$ and σ_k is *not* a disjunction, put $\Gamma_{k+1} = \Gamma_k \cup \{\sigma_k\}$.
3. If $\Gamma_k \cup \{\sigma_k\} \not\vdash \gamma$ and σ_k is a disjunction $\alpha \vee \beta$, then (a) $\Gamma_k \cup \{\sigma_k, \alpha\} \not\vdash \gamma$ or (b) $\Gamma_k \cup \{\sigma_k, \beta\} \not\vdash \gamma$. If (a) holds, put $\Gamma_{k+1} = \Gamma_k \cup \{\sigma_k, \alpha\}$; if (b), put $\Gamma_{k+1} = \Gamma_k \cup \{\sigma_k, \beta\}$.

Now define $\Pi = \bigcup_k \Gamma_k$. It follows immediately from 1.–3. that $\Gamma_k \not\vdash \gamma \Rightarrow \Gamma_{k+1} \not\vdash \gamma$, so that $\Gamma_k \not\vdash \gamma$ for all k , whence $\Pi \not\vdash \gamma$. Moreover, Π is a theory. For if $\Pi \vdash \sigma_k$, then since $\Pi \not\vdash \gamma$, $\Pi \cup \{\sigma_k\} \not\vdash \gamma$, so $\Gamma_k \cup \{\sigma_k\} \not\vdash \gamma$, whence $\sigma_k \in \Gamma_{k+1} \subseteq \Pi$.

And finally, Π is prime. For if $\alpha \vee \beta \in \Pi$ with $\alpha \vee \beta = \sigma_k$, then $\Pi \cup \{\sigma_k\} \vdash \gamma$, so that $\Gamma_k \cup \{\sigma_k\} \vdash \gamma$, whence $\Gamma_{k+1} = \Gamma_k \cup \{\sigma_k, \alpha\}$ or $\Gamma_{k+1} = \Gamma_k \cup \{\sigma_k, \beta\}$. Therefore $\alpha \in \Gamma_{k+1} \subseteq \Pi$ or $\beta \in \Gamma_{k+1} \subseteq \Pi$. ■

Given a consistent set of sentences Γ , we define the *canonical frame* associated with Γ to be the frame $\mathbb{K}_\Gamma = (P_\Gamma, \subseteq, \Sigma_\Gamma)$, where P_Γ is the set of prime theories extending Γ , and, for $\Delta \in P_\Gamma$, $\Sigma_\Gamma(\Delta)$ is the set of atomic sentences in Δ . For this frame we have the

Fundamental Lemma. (1) For all $\Delta \in P_\Gamma$, all α , $\Delta \Vdash_{\mathbb{K}_\Gamma} \alpha \Leftrightarrow \alpha \in \Delta$.

(2) $\Vdash_{\mathbb{K}_\Gamma} \alpha \Leftrightarrow \Gamma \vdash \alpha$; in particular $\Vdash_{\mathbb{K}_\Gamma} \Gamma$.

Proof. (1) is proved by induction on the number of logical symbols in α . For α atomic it holds by the definition of Σ_Γ . The induction step for \wedge is trivial and that for \vee follows immediately from the primeness of Δ . To establish the induction step for \rightarrow , we argue as follows. Supposing that (1) holds for α and β , we have:

$$\begin{aligned} \Delta \Vdash_{\mathbb{K}_\Gamma} \alpha \rightarrow \beta &\Leftrightarrow \forall \Delta' \supseteq \Delta. \Delta' \Vdash_{\mathbb{K}_\Gamma} \alpha \Rightarrow \Delta' \Vdash_{\mathbb{K}_\Gamma} \beta \\ &\Leftrightarrow \forall \Delta' \supseteq \Delta. \alpha \in \Delta' \Rightarrow \beta \in \Delta' \\ &\Leftrightarrow^* \forall \Delta' \supseteq \Delta. (\alpha \rightarrow \beta) \in \Delta' \\ &\Leftrightarrow \alpha \rightarrow \beta \in \Delta. \end{aligned}$$

We need to justify the equivalence marked *: clearly $\alpha \rightarrow \beta \in \Delta' \Rightarrow [\alpha \in \Delta' \Rightarrow \beta \in \Delta']$. Conversely suppose $\alpha \rightarrow \beta \notin \Delta'$ for some $\Delta' \supseteq \Delta$. Then $\Delta' \cup \{\alpha\} \vdash \beta$, so by the extension lemma there is $\Delta'' \in P_\Gamma$

such that $\beta \notin \Delta''$ and $\Delta' \cup \{\alpha\} \subseteq \Delta''$. Hence $\alpha \in \Delta'' \Rightarrow \beta \in \Delta''$.

Thus (1) is proved.

(2). Clearly $\Gamma \vdash \alpha \Rightarrow \alpha \in \Delta$ for all $\Delta \in P_\Gamma \Rightarrow \Vdash_{K_\Gamma} \alpha$ by (1).

Conversely if $\Gamma \not\vdash \alpha$ there is $\Delta \in P_\Gamma$ with $\alpha \notin \Delta$. Then $\Delta \not\vdash_{K_\Gamma} \alpha$ by (1), whence $\not\vdash_{K_\Gamma} \alpha$. ■

All this leads to the

Completeness Theorem. $\Gamma \Vdash \alpha \Rightarrow \Gamma \vdash \alpha$.

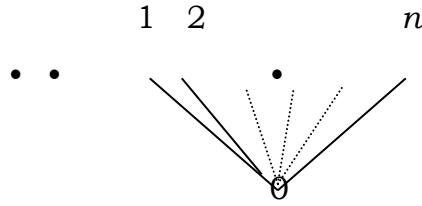
Proof. If $\Gamma \Vdash \alpha$ then since $\Vdash_{K_\Gamma} \Gamma$ it follows that $\Vdash_{K_\Gamma} \alpha$, whence $\Gamma \vdash \alpha$. ■

THE DISJUNCTION PROPERTY

Kripke semantics can be used to establish other significant facts about intuitionistic logic. For example, in 1933 Gödel proved that no finite truth-table fully characterizes intuitionistic propositional logic. This is easily proved using frames. For if n -valued truth tables characterized such logic, then, under any assignment of truth values, of any $n + 1$ atomic sentences p_0, p_1, \dots, p_n , at least two would obtain the same value. Accordingly, the sentence

$$\sigma = \bigvee_{0 \leq i < j \leq n} p_i \leftrightarrow p_j$$

would have to be true in all frames. However, consider the following frame:



Here $S(0) = \emptyset$ and, for each i , $1 \leq i \leq n$, $S(i) = \{p_i\}$. In this frame, clearly $0 \not\Vdash \sigma$.

Both propositional and first-order intuitionistic logic possess the important *disjunction property*: for sentences α , β , if $\vdash \alpha \vee \beta$, then $\vdash \alpha$ or $\vdash \beta$. Using frames, we prove this in the propositional case.

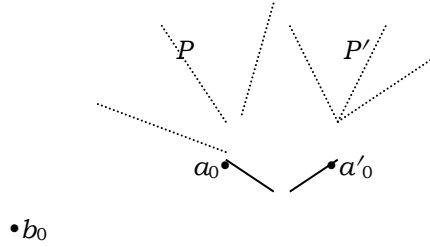
First, some definitions. A *bottom element* of a partially ordered set P is an element $a_0 \in P$ such that $a_0 \leq a$ for all $a \in P$. A bottom element of a partially ordered set is also referred to as a bottom element of any frame built on it. A subset Q of a P is said to be *closed* if $a \in Q$, $a \leq b \Rightarrow b \in Q$. Given a frame $\mathbb{K} = (P, \leq, S)$ built on P , any closed subset Q of P determines a frame $\mathbb{K}|Q = (Q, \leq, S')$ —called the *restriction* of \mathbb{K} to Q —with $S'(a) = S(a)$ for $a \in Q$. It is easily proved by induction on sentences α that, for any $a \in Q$,

$$a \Vdash_{\mathbb{K}} \alpha \Leftrightarrow a \Vdash_{\mathbb{K}|Q} \alpha.$$

We next show that, if Γ is a set of sentences and γ a sentence such that $\Gamma \not\vdash \gamma$, there is a frame \mathbb{K} with a bottom element a_0 such that $a_0 \Vdash_{\mathbb{K}} \Gamma$ and $a_0 \not\Vdash_{\mathbb{K}} \gamma$. To prove this, let Π_0 be a prime theory

extending Γ such that $\gamma \notin \Pi_0$ and let K be the restriction of K_Γ to the closed subset $\{\Delta: \Pi_0 \subseteq \Delta\}$ of P_Γ . Then K has bottom element Π_0 ; the fundamental lemma implies $\Pi_0 \Vdash_{K_\Gamma} \Gamma$ and $\Pi_0 \not\Vdash_{K_\Gamma} \gamma$; and it follows from this and the fact above that $\Pi_0 \Vdash_K \Gamma$ and $\Pi_0 \not\Vdash_K \gamma$.

Now we can show that intuitionist propositional logic has the disjunction property. For suppose that both $\not\Vdash \alpha$ and $\not\Vdash \beta$. Then by the above there are frames $K = (P, \leq, S)$ and $K' = (P', \leq', S')$ with bottom elements a_0, a'_0 for which $a_0 \not\Vdash_K \alpha$ and $a'_0 \not\Vdash_{K'} \beta$. Without loss of generality we may, and do, assume that P and P' are disjoint. Let $Q = P \cup P' \cup \{b_0\}$, where b_0 is some element outside $P \cup P'$, and let \triangleleft be the partial order on Q with bottom element b_0 which coincides with \leq on P and with \leq' on P' . Clearly P and P' are then closed subsets of Q .



Let $Q = (Q, \triangleleft, T)$ be the frame with $T(b_0) = \emptyset$, $T(a) = S(a)$ for $a \in P$, $T(a') = S'(a')$ for $a' \in P'$. Then for $a \in P$, $a' \in P'$, we have

$$a \Vdash_K \alpha \Leftrightarrow a \Vdash_Q \alpha \quad a' \Vdash_{K'} \beta \Leftrightarrow a' \Vdash_Q \beta.$$

So $a_0 \Vdash_Q \alpha$, $a_0' \Vdash_Q \beta$, whence $b_0 \Vdash_Q \alpha$, $b_0 \Vdash_Q \beta$, and $b_0 \Vdash_Q \alpha \vee \beta$. We conclude that

$\Vdash_Q \alpha \vee \beta$, and $\not\vdash \alpha \vee \beta$ follows by soundness. This establishes the disjunction property.

It can be shown that, in addition to possessing the disjunction property, intuitionistic predicate logic¹³ has the *existence property*: if $\vdash \exists x\alpha(x)$, then $\vdash \alpha(t)$ for some closed term t .

INTUITIONISTIC LOGIC IN LINEAR STYLE

Intuitionistic logic can also be presented in traditional linear style. We now suppose that \mathcal{L} has the equality symbol $=$. The system of *intuitionistic first-order logic* in \mathcal{L} has the following *axioms* and *rules of inference*:

Axioms

$$\alpha \rightarrow (\beta \rightarrow \alpha) \quad [\alpha \rightarrow (\beta \rightarrow \gamma) \rightarrow [(\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \gamma)]]$$

$$\begin{array}{l} \alpha \rightarrow (\beta \rightarrow \alpha \wedge \beta) \quad \alpha \wedge \beta \rightarrow \alpha \quad \alpha \wedge \beta \rightarrow \beta \\ \alpha \rightarrow \alpha \vee \beta \quad \beta \rightarrow \alpha \vee \beta \quad (\alpha \rightarrow \gamma) \rightarrow [(\beta \rightarrow \gamma) \rightarrow (\alpha \vee \beta \rightarrow \gamma)] \\ [\alpha \rightarrow (\beta \rightarrow \gamma) \rightarrow [(\alpha \rightarrow \beta) \rightarrow (\alpha \rightarrow \gamma)]] \end{array}$$

$$(\alpha \rightarrow \beta) \rightarrow [(\alpha \rightarrow \neg\beta) \rightarrow \neg\alpha] \quad \neg\alpha \rightarrow (\alpha \rightarrow \beta)$$

$$\alpha(t) \rightarrow \exists x\alpha(x) \quad \forall x\alpha(x) \rightarrow \alpha(y) \quad (x \text{ free in } \alpha \text{ and } t \text{ free for } x \text{ in } \alpha)$$

$$x = x \quad x = y \rightarrow y = x \quad \alpha(x) \wedge x = y \rightarrow \alpha(y) \quad (x \text{ free for } y \text{ in } \alpha)$$

¹³ with no function symbols and at least one constant symbol.

Rules of Inference

$$\textit{Modus ponens} \quad \frac{\alpha, \alpha \rightarrow \beta}{\beta}$$

$$\textit{Quantifier rules} \quad \frac{\beta \rightarrow \alpha(x)}{\beta \rightarrow \forall x \alpha(x)} \quad \frac{\alpha(x) \rightarrow \beta}{\exists x \alpha(x) \rightarrow \beta}$$

(x not free in β)

In each of the rules of inference the formula below the line is called an *immediate consequence* of the formula(s) above the line.

The system of *free* first-order intuitionistic logic is obtained by restricting the modus ponens rule to cases where all variables free in α are also free in β . This allows for the possibility of empty domains of interpretation.

If Γ is a set of formulas, and α a formula, of \mathcal{L} , a(n) (intuitionistic) *proof* of α from Γ is a sequence $\alpha_1, \dots, \alpha_n$ of formulas such that α_n is α and, for any j , $1 \leq j \leq n$, α_j is either an axiom, a member of Γ , or is an immediate consequence of some α_k with $k < j$. If there exists a proof of α from Γ , we write $\Gamma \vdash \alpha$ and say that α is *provable* from Γ . α is a *theorem* of intuitionistic logic, written $\vdash \alpha$, if $\emptyset \vdash \alpha$.

If to the axioms above we add the law of excluded middle $\alpha \vee \neg \alpha$ or the law of double negation $\neg \neg \alpha \rightarrow \alpha$, then we obtain classical first-order logic.

HEYTING ALGEBRAS AND ALGEBRAIC INTERPRETATIONS OF
INTUITIONISTIC LOGIC

We now introduce the idea of an *algebraic interpretation* of intuitionistic logic. To do this we require the concept of a lattice.

A *lattice* is a partially ordered set L with partial ordering \leq in which each two-element subset $\{x, y\}$ has a supremum or *join*—denoted by $x \vee y$ —and an infimum or *meet*—denoted by $x \wedge y$. A lattice L is *complete* if every subset X (including \emptyset) has a supremum or *join*—denoted by $\bigvee X$ —and an infimum or *meet*—denoted by $\bigwedge X$. Note that $\bigvee \emptyset = 0$, the least or *bottom* element of L , and $\bigwedge \emptyset = 1$, the largest or *top* element of L .

A *Heyting algebra* is a lattice L with top and bottom elements such that, for any elements $x, y \in L$, there is an element—denoted by $x \Rightarrow y$ —of L such that, for any $z \in L$,

$$z \leq x \Rightarrow y \text{ iff } z \wedge x \leq y.$$

Thus $x \Rightarrow y$ is the *largest* element z such that $z \wedge x \leq y$. So in particular, if we write x^* for $x \Rightarrow 0$, then x^* is the largest element z such that $x \wedge z = 0$: it is called the *pseudocomplement* of x .

A *Boolean algebra* is a Heyting algebra in which $x^{**} = x$ for all x , or equivalently, in which $x \vee x^* = 1$ for all x .

Heyting algebras are related to intuitionistic propositional logic in precisely the same way as Boolean algebras are related to

classical propositional logic. That is, suppose given a propositional language; let \mathcal{P} be its set of propositional variables. Given a map $f: \mathcal{P} \rightarrow L$ to a Heyting algebra L , we extend f to a map $\alpha \mapsto \llbracket \alpha \rrbracket$ of the set of formulas of \mathcal{L} to L à la Tarski:

$$\llbracket \alpha \wedge \beta \rrbracket = \llbracket \alpha \rrbracket \wedge \llbracket \beta \rrbracket \quad \llbracket \alpha \vee \beta \rrbracket = \llbracket \alpha \rrbracket \vee \llbracket \beta \rrbracket \quad \llbracket \alpha \Rightarrow \beta \rrbracket = \llbracket \alpha \rrbracket \Rightarrow \llbracket \beta \rrbracket$$

$$\llbracket \neg \alpha \rrbracket = \llbracket \alpha \rrbracket^*$$

A formula α is said to be (Heyting) *valid*—written $\vdash \alpha$ —if $\llbracket \alpha \rrbracket = 1$ for any such map f . It can then be shown that α is valid iff $\vdash \alpha$ in the intuitionistic propositional calculus, i.e., iff α is provable from the propositional axioms listed above.

A basic fact about *complete* Heyting algebras is that the following identity holds in them:

$$(*) \quad x \wedge \bigvee_{i \in I} \dots \bigvee_{i \in I} \dots$$

And conversely, in any complete lattice satisfying (*), defining the operation \Rightarrow by $x \Rightarrow y = \bigvee \{z: z \wedge x \leq y\}$ turns it into a Heyting algebra.

To prove this, we observe that in any complete Heyting algebra,

$$\begin{aligned}
& x \wedge \bigvee_{i \in I} y_i \leq \bigvee_{i \in I} (y_i \wedge x) ; \\
& \leftrightarrow y_i \leq x \Rightarrow z, \text{ all } i \\
& \leftrightarrow y_i \wedge x \leq z, \text{ all } i \\
& \leftrightarrow \bigvee_{i \in I} (y_i \wedge x) \leq z
\end{aligned}$$

Conversely, if (*) is satisfied and $x \Rightarrow y$ is defined as above, then

$$(x \Rightarrow y) \wedge x \leq \bigvee \{z: z \wedge x \leq y\} \wedge x = \bigvee \{z \wedge x: z \wedge x \leq y\} \leq y$$

So $z \leq x \Rightarrow y \rightarrow z \wedge x \leq (x \Rightarrow y) \wedge x \leq y$. The reverse inequality is an immediate consequence of the definition.

In view of this result a complete Heyting algebra is frequently defined to be a complete lattice satisfying (*).

Complete Heyting algebras are related to intuitionistic first-order logic in the same way as complete Boolean algebras are related to classical first-order logic. To be precise, let \mathcal{L} be a first-order language whose sole extralogical symbol is a binary predicate symbol P . An \mathcal{L} -structure is a quadruple $\mathbf{M} = (M, eq, Q, L)$, where M is a nonempty set, L is a complete Heyting algebra and eq and Q are maps $M^2 \rightarrow M$ satisfying, for all $m, n, m', n' \in M$,

$$eq(m, m) = 1, \quad eq(m, n) = eq(n, m), \quad eq(m, n) \wedge eq(n, n') \leq eq(m, n'),$$

$$Q(m, n) \wedge eq(m, m') \leq Q(m', n), \quad Q(m, n) \wedge eq(n, n') \leq Q(m, n').$$

For any formula α of \mathcal{L} and any finite sequence $\mathbf{x} = \langle x_1, \dots, x_n \rangle$ of variables of \mathcal{L} containing all the free variables of α , we define for any \mathcal{L} -structure \mathbf{M} a map

$$\llbracket \alpha \rrbracket^{\mathbf{M}_x}: M^n \rightarrow L$$

recursively as follows:

$$\llbracket x_p = x_q \rrbracket^{\mathbf{M}_x} = \langle m_1 \dots, m_n \rangle \mapsto eq(m_p, m_q),$$

$$\llbracket Px_p x_q \rrbracket^{\mathbf{M}_x} = \langle m_1 \dots, m_n \rangle \mapsto Q(m_p, m_q),$$

$$\llbracket \alpha \wedge \beta \rrbracket^{\mathbf{M}_x} = \llbracket \alpha \rrbracket^{\mathbf{M}_x} \wedge \llbracket \beta \rrbracket^{\mathbf{M}_x}, \text{ and similar clauses for the other}$$

connectives,

$$\llbracket \exists y \alpha \rrbracket^{\mathbf{M}_x} = \langle m_1 \dots, m_n \rangle \mapsto \bigvee_{m \in M} \llbracket \alpha(y/u) \rrbracket^{\mathbf{M}_{ux}(m, m_1 \dots, m_n)}$$

$$\llbracket \forall y \alpha \rrbracket^{\mathbf{M}_x} = \langle m_1 \dots, m_n \rangle \mapsto \bigwedge_{m \in M} \llbracket \alpha(y/u) \rrbracket^{\mathbf{M}_{ux}(m, m_1 \dots, m_n)}$$

Call α **M**-*valid* if $\llbracket \alpha \rrbracket^{\mathbf{M}_x}$ is identically 1, where \mathbf{x} is the sequence of all free variables of α . Then it can be shown that α is **M**-*valid for all M* iff α is provable in intuitionistic first-order logic. This is the *algebraic completeness theorem* for intuitionistic first-order logic. A similar result may be obtained for free intuitionistic logic by allowing the domains of \mathcal{L} -structures to be empty.

INTUITIONISTIC FIRST-ORDER ARITHMETIC

Finally, we make some observations on the first-order intuitionistic theory of the natural numbers.

Heyting or *intuitionistic arithmetic* **HA** is formulated within the first-order *language of arithmetic*, which has symbols $+$, \cdot , s , 0 , 1 . The axioms of **HA** are the usual ones, viz.,

1. $sx = sy \rightarrow x = y$
2. $\neg sx = 0$
3. $x + 0 = x \quad x + sy = s(x + y)$
4. $x \cdot 0 = 0 \quad x \cdot sy = x \cdot y + x$
5. $\alpha(0) \wedge \forall x(\alpha(x) \rightarrow \alpha(sx)) \rightarrow \forall x \alpha(x)$.

Axiom 5 is the *principle of mathematical induction*. Using this, one can establish the decidability of the equality relation:

$$\mathbf{HA} \vdash \forall x \forall y (x = y \vee x \neq y).$$

The ordering relations $<$ and \leq are defined by $x < y \Leftrightarrow \exists z (y = x + sz)$ and $x \leq y \Leftrightarrow x < y \vee x = y$. Using induction one can prove the *trichotomy principle*:

$$\mathbf{HA} \vdash \forall x \forall y (x < y \vee x = y \vee y < x).$$

In classical arithmetic as an immediate consequence of the principle of induction one obtains the *least number principle*, viz.,

$$\exists x \alpha(x) \rightarrow \exists x [\alpha(x) \wedge \forall y (\alpha(y) \rightarrow x \leq y)].$$

In Heyting arithmetic, however, this principle cannot be derived¹⁴, since, as the following simple argument shows, it implies the law of excluded middle. Let β be any sentence and let $\alpha(x)$ be the formula $\beta \vee x \neq 0$. Then clearly $\exists x \alpha(x)$, so if the least number principle held there would exist n_0 for which $\alpha(n_0)$ and $\forall y (\alpha(y) \rightarrow n_0 \leq y)$, that is,

¹⁴ But it can be derived from the assumption that α satisfies the law of excluded middle, i.e. $\forall x (\alpha(x) \vee \neg \alpha(x))$.

$$(1) \beta \vee n_0 \neq 0 \quad (2) \forall y(\beta \vee y \neq 0 \rightarrow n_0 \leq y).$$

From (1) it follows that $n_0 = 0 \rightarrow \beta$, and from (2) that $\beta \rightarrow n_0 = 0$. Therefore $n_0 = 0 \leftrightarrow \beta$, whence $n_0 \neq 0 \rightarrow \neg\beta$. Since $\mathbf{HA} \vdash n_0 = 0 \vee n_0 \neq 0$, we infer $\beta \vee \neg\beta$.

\mathbf{HA} also has the disjunction and existence properties: in fact, if $\mathbf{HA} \vdash \exists x \alpha(x)$, then $\mathbf{HA} \vdash \alpha(\mathbf{n})$ for some n , where \mathbf{n} is the closed term $s\dots s0$ with n s 's.

4. Interlude: Constructivity in Mathematics before Brouwer

Nonconstructive proofs in mathematics are an essentially modern conception: with singularly few exceptions, all mathematical proofs before 1880 were constructive. Indeed, the very notion of “existence” in mathematics was, to all intents and purposes, taken to mean “constructive existence”.

There were, however, a few nonconstructive proofs, for example, Euler’s proof in the 18th century of the existence of infinitely many prime numbers from his formula

$$\prod_{p \text{ prime}} (1 - p^{-s})^{-1} = \sum_{n=1}^{\infty} n^{-s} :$$

if there were only finitely many primes p , the product would converge for $s = 1$, but the sum is known to diverge. (Of course, the existence of infinitely many primes is constructively provable.) Another example, already mentioned, is the proof of the fundamental theorem of algebra using Liouville’s theorem, but again, this has a constructive proof. Hilbert became celebrated for his nonconstructive proof of the finite basis theorem for polynomial ideals, causing his colleague Gordan to exclaim “this is not mathematics, it is theology!” Hilbert also supplied an entirely nonconstructive proof of Waring’s conjecture that, for each number m , there is a number n such that every number is the sum of not more than m n th powers.

But it was Cantor’s development of set theory, with its embrace of the actual infinite, which truly opened the door to the unrestricted use of nonconstructive arguments in mathematics. This provoked some reaction, especially from the German

mathematician Kronecker, the most prominent of Cantor's intellectual opponents, who observed in 1886 that

God made the natural numbers, everything else is the work of Man.

Kronecker also rejected the notion of an arbitrary sequence of natural numbers, asserting in 1889:

Even the general concept of an infinite series, for example, one in which only specified powers appear, is in my opinion only permissible with the condition that in each special case, on the basis of the arithmetical formation laws of the coefficients, certain hypotheses are satisfied which permit one to reduce the series to a finite expression—which thus actually makes the extension of the concept of a finite sequence unnecessary.

The issue came to a head in 1904 with the publication of Zermelo's proof of the well-ordering theorem that any set can be ordered in such a way as to ensure that every nonempty subset has a least element. In his proof Zermelo had formulated and made essential use of the *axiom of choice*, which asserts that, given any family of nonempty sets A , there is a function—a *choice function*— f defined on A such that $f(A) \in A$ for each $A \in A$. The “nonconstructive” character of this principle provoked the objections of a number of prominent mathematicians of the day. Borel, for example, claimed that what Zermelo had actually done was to demonstrate the equivalence of the problems of (1) well-ordering an arbitrary set M and (2) choosing a distinguished element from each nonempty subset of M . What Zermelo had *failed*

to show, according to Borel, was that the equivalence of (1) and (2) furnishes

a general solution to the first problem. In fact, to regard the second problem as resolved for a given set M , one needs a means , at least a theoretical one, for determining a distinguished element m' from an arbitrary subset M' of M ; and this problem appears to be one of the most difficult, if one supposes, for the sake of definiteness, that coincides with the continuum..

In using the word “determining” here Borel is evidently demanding that the selection of a distinguished element from an arbitrary subset of a set be made *constructively*. This requirement is left completely unaddressed by the axiom of choice. Having come to regard the idea of an uncountable set as fundamentally vague, he was particularly unhappy with Zermelo’s use of the axiom of choice to make uncountably many arbitrary “choices”, as was required when establishing the well-orderability of the continuum.

The French mathematician Baire’s objections went still further. Like Kronecker, he rejected the completed infinite altogether, and even regarded the potential infinite as a mere *façon de parler*. He went so far as to assert that, even were one to be given an infinite set,

I consider it false to regard the subsets of this set as being given.

For Baire, in the last analysis, everything in mathematics must be reduced to the finite.

Lebesgue put the central question in unequivocally constructive terms: *Can the existence of a mathematical object be proved without at the same time defining it?* Lebesgue says, in essence, no—thus bringing him into the constructivist camp. He rejected proofs that demonstrate the existence of a nonempty class of objects of a certain kind as opposed to actually producing an object of that kind. He also objected to the idea of making an infinity, even a countable infinity, of arbitrary choices.

Among classical mathematicians, the term “constructive” is still sometimes used with the meaning “without making use of the axiom of choice”.

5. Intuitionistic Set Theory.

INTUITIONISTIC ZERMELO SET THEORY

The system \mathbf{Z}_I of *intuitionistic Zermelo set theory* is formulated in the usual first-order language of set theory with relation symbols $=, \in$ but is subject to the axioms and rules of intuitionistic first-order logic. Arguments in \mathbf{Z}_I will be presented informally; in particular we shall make use of the standard notations of classical set theory: $\exists y \in x, \forall y \in x, \{x: \alpha\}, x \cup y, Px, (x, y), x \subseteq y, \emptyset, 0, 1, 2$, etc. The *axioms* of \mathbf{Z}_I are *Extensionality, Pairing, Union, Power set, Infinity* and *Separation*:

$$\mathbf{Ext} \quad \forall x \forall y [\forall z (z \in x \leftrightarrow z \in y) \leftrightarrow x = y]$$

$$\mathbf{Pair} \quad \forall x \forall y \exists z \forall w (w \in z \leftrightarrow w = x \vee w = y)$$

$$\mathbf{Union} \quad \forall x \exists z \forall w (w \in z \leftrightarrow \exists y \in x. w \in y)$$

$$\mathbf{Power} \quad \forall x \exists z \forall w (w \in z \leftrightarrow w \subseteq x)$$

$$\mathbf{Inf} \quad \exists x (\emptyset \in x \wedge \forall y \in x. y \cup \{y\} \in x)$$

$$\mathbf{Sep} \quad \exists z \forall w (w \in z \leftrightarrow w \in x \wedge \alpha).$$

For any set A , PA is a complete Heyting algebra with operations \cup, \cap and \Rightarrow , where $U \Rightarrow V = \{x: x \in U \rightarrow x \in V\}$, and top and bottom element A and \emptyset respectively.

We write $\{\tau | \alpha\}$ for $\{x: x = \tau \wedge \alpha\}$ where τ is a closed term: without the law of excluded middle we cannot conclude that $\{\tau | \alpha\} = \emptyset$ or $\{\tau\}$. From **Ext** we infer that $\{\tau | \alpha\} = \{\tau | \beta\} \Leftrightarrow (\alpha \leftrightarrow \beta)$; thus, in particular, the elements of $P1$ correspond naturally to *truth values*, i.e. propositions identified under equivalence. $P1$ is called the (Heyting) *algebra of truth values* and is denoted by Ω . The top

element 1 of Ω is usually written *true* and the bottom element 0 as *false*.

Properties of Ω correspond to logical properties of the set theory. Thus, for instance, the law of excluded middle $\alpha \vee \neg \alpha$ and the weak law of excluded middle $\neg \alpha \vee \neg \neg \alpha$ (equivalent to de Morgan's law $\neg(\alpha \wedge \beta) \rightarrow \neg \alpha \vee \neg \beta$) correspond respectively to the properties

LEM $\forall \omega \in \Omega. \omega = \text{true} \vee \omega = \text{false}$

WLEM $\forall \omega \in \Omega. \omega = \text{false} \vee \omega \neq \text{false}$.

Calling a set A *decidable* if $\forall x \in A \forall y \in A. x = y \vee x \neq y$, each of the following is equivalent to **LEM**:

1. *Every set is decidable*
2. *Ω is decidable*
3. *Membership is decidable: $\forall x \forall y (x \in y \vee x \notin y)$*
4. $\forall x (0 \in x \vee 0 \notin x)$
5. *$(2, \leq)$ is well-ordered.*

(To show that 5. implies **LEM**, observe that the least element of $\{0 \mid \alpha\} \cup \{1\} \subseteq 2$ is either 0 or 1; if it is 0, α must hold, and if it is 1, α must fail.)

Using the axiom of infinity, the set \mathbb{N} of natural numbers can be constructed as usual. \mathbb{N} is decidable and satisfies the familiar Peano axioms including induction, but it is well-ordered only if **LEM** holds. In fact **LEM** also follows from the *domino principle* for \mathbb{N} :

$$\alpha(0) \wedge \exists n \neg \alpha(n) \rightarrow \exists n [\alpha(n) \wedge \neg \alpha(n+1)]^{15}.$$

To see this, take any proposition β and define $\alpha(n)$ to be the formula $n = 0 \vee (n = 1 \wedge \beta)$. Then clearly $\alpha(0) \wedge \exists n \neg \alpha(n)$ holds, so we infer from the domino principle that there is n_0 for which $\alpha(n)$ and $\neg \alpha(n+1)$, i.e.,

$$(*) \quad n_0 = 0 \vee (n_0 = 1 \wedge \beta)$$

and

$$\neg(n_0 + 1 = 1 \wedge \beta)$$

whence

$$\neg(n_0 = 0 \wedge \beta).$$

From this last we infer $n_0 = 0 \rightarrow \neg\beta$, which, together with (*), gives $\beta \vee \neg\beta$.

The notion of a *function* is defined as usual in **ZF_I**; we employ the standard notations for functions. A *choice function* on a set A is a function f with domain A such that $f(a) \in a$ whenever $\exists x.x \in a$. The *axiom of choice* **AC** is the assertion that every set has a choice function. Remarkably, **AC** implies **LEM**; in fact we have the

Theorem. If each doubleton has a choice function, then **LEM** holds (and conversely).

¹⁵ Here and in the sequel we shall use lower case letters m, n as variables ranging over \mathbb{N} .

Proof. Define $U = \{x \in 2: x = 0 \vee \alpha\}$ and $V = \{x \in 2: x = 1 \vee \alpha\}$, and suppose given a choice function f on $\{U, V\}$. Writing $a = f(U)$, $b = f(V)$, we then have $a \in U, b \in V$, i.e.

$$(a = 0 \vee \alpha) \wedge (b = 1 \vee \alpha).$$

Hence

$$a = 0 \wedge (b = 1 \vee \alpha),$$

whence

$$a \neq b \vee \alpha. \quad (*)$$

But

$$\alpha \rightarrow U = V \rightarrow a = b,$$

so that

$$a \neq b \rightarrow \neg\alpha.$$

This, together with (*), gives $\alpha \vee \neg\alpha$. ■

It can also be shown that the assertion *any singleton has a choice function* is equivalent in \mathbf{Z}_1 to the (constructively invalid) “independence of premises” rule,

$$\frac{\alpha \rightarrow \exists x (x \in A \wedge \beta(x))}{\exists x (\alpha \rightarrow x \in A \wedge \beta(x))}.$$

In classical set theory one proves the well-known *Schröder-Bernstein theorem*: if each of two sets A and B can be injected into the other, then there is a bijection between A and B . This is usually derived as a consequence of the proposition

SB: for any set X and any injection $f: X \rightarrow X$ there is a bijection $h: X \rightarrow X$ such that $h \subseteq f \cup f^{-1}$, i.e., $\forall x \in X. h(x) = f(x) \vee f(h(x)) = x$.

In **Z_I** this assertion implies (and so is equivalent to) **LEM**. Here is the proof.

Define, for any proposition α ,

$$\mathbb{N}^\alpha = \mathbb{N} - \{0\} \cup \{0 \mid \alpha\} \quad f = \{(n, n+1) : n \neq 0\} \cup \{(0, 1) \mid \alpha\}.$$

Then $f: \mathbb{N}^\alpha \rightarrow \mathbb{N}^\alpha$. Clearly

$$(*) \quad 1 \in \text{range}(f) \leftrightarrow 0 \in \mathbb{N}^\alpha \leftrightarrow \alpha.$$

Now suppose given a bijection $h: \mathbb{N}^\alpha \rightarrow \mathbb{N}^\alpha$ such that

$$\forall x \in \mathbb{N}^\alpha. h(x) = f(x) \vee f(h(x)) = x.$$

If α holds, then f is just the usual successor function on \mathbb{N} ($= \mathbb{N}^\alpha$) and so

$$\alpha \wedge h(n) = 0 \rightarrow h(n) \neq f(n) \rightarrow 1 = f(0) = f(h(n)) = n \rightarrow n = 1,$$

whence

$$\alpha \rightarrow h(1) = 0.$$

Thus

$$(**) \quad h(1) \neq 0 \rightarrow \neg\alpha$$

But

$$h(1) = f(1) \vee f(h(1)) = 1.$$

The first disjunct implies $h(1) \neq 0$ and (**) gives $\neg\alpha$. From the second disjunct we infer $1 \in \text{range}(f)$ and (*) yields α . Thus we have derived $\alpha \vee \neg\alpha$.

In classical set theory Zorn's lemma¹⁶ is used to prove the so-called *order extension principle*, namely: every partial ordering on a set can be extended to a total ordering. We will show that this principle implies the intuitionistically invalid law $\alpha \rightarrow \beta \vee \beta \rightarrow \alpha$.

To prove this, we first observe that if $U, V \subseteq 1$, then

$$(*) \quad (U = 1 \rightarrow V = 1) \leftrightarrow U \subseteq V.$$

Now suppose that \leq is a partial order on Ω extending \subseteq . Then $U \leq 1$ for all $U \subseteq 1$. Now

$$U \leq V \wedge U = 1 \rightarrow 1 \leq V \rightarrow V = 1,$$

whence, using (*),

$$U \leq V \rightarrow (U = 1 \rightarrow V = 1) \rightarrow U \subseteq V.$$

We conclude that \leq and \subseteq coincide. Accordingly, if \subseteq could be extended to a total order on Ω , \subseteq would have to be a total order on

¹⁶ Zorn's lemma, although classically equivalent to the axiom of choice, is not intuitionistically equivalent to it. In fact it can be shown that, unlike the axiom of choice, which implies **LEM**, Zorn's lemma has no nonconstructive consequences whatsoever.

Ω itself. But this is clearly tantamount to the truth of $\alpha \rightarrow \beta \vee \beta \rightarrow \alpha$ for arbitrary propositions α and β .

The negation operation \neg on propositions corresponds to the complementation operation on Ω ; we use the same symbol \neg to denote the latter. This operation of course satisfies

$$\omega \subseteq \neg\omega' \leftrightarrow \omega \cap \omega' = \text{false}.$$

Classically, \neg also satisfies the dual law, viz.

$$\neg\omega \subseteq \omega' \leftrightarrow \omega \cup \omega' = \text{true}.$$

But intuitionistically, this is far from being the case. Indeed, the assumption that there exists *any* operation $\neg: \Omega \rightarrow \Omega$ satisfying

$$\neg\omega \subseteq \omega' \leftrightarrow \omega \cup \omega' = \text{true}$$

implies (and so is equivalent to) **LEM**. For suppose such an operation existed. Then

$$\neg\text{true} \subseteq \text{false} \leftrightarrow \text{false} \cup \text{true} = \text{true},$$

so that $\neg\text{true} \subseteq \text{false}$, whence $\neg\text{true} = \text{false}$. Next,

$$0 \in \neg\omega \wedge 0 \in \omega \rightarrow 0 \in \neg\omega \wedge \omega = \text{true} \rightarrow 0 \in \neg\text{true} = \text{false}.$$

Since $0 \notin \text{false}$, it follows that

$$0 \in \neg\omega \rightarrow 0 \notin \omega \rightarrow 0 \in \neg\omega,$$

and from this we infer that $\neg\omega \subseteq \neg\omega$. Since, obviously, $\omega \cup \neg\omega = \text{true}$, it then follows that, for any ω , $\omega \cup \neg\omega = \text{true}$, which is **LEM**.

DEFINITIONS OF FINITENESS.

Fix a set E ; by “set”, “family” etc. we shall for the time being mean “subset of E ”, “family of subsets of E ”, etc.

A family F is

- (a) *strictly inductive* if $\emptyset \in F \wedge \forall X \in F \forall x \in E - X. X \cup \{x\} \in F$.
- (b) *inductive* if $\emptyset \in F \wedge \forall X \in F \forall x \in E. X \cup \{x\} \in F$.
- (c) *K(uratowski)-inductive* if $\emptyset \in F \wedge \forall x \in E . \{x\} \in F \wedge \forall XY \in F. X \cup Y \in F$.

The members of the least $\left\{ \begin{array}{l} \text{strictly inductive} \\ \text{inductive} \\ \text{K-inductive} \end{array} \right\}$ families

are called

$\left\{ \begin{array}{l} \text{strictly finite} \\ \text{finite} \\ \text{K-finite.} \end{array} \right.$

It can be shown that $\mathbf{Z}_I \vdash \text{strictly finite} \rightarrow \text{finite} \leftrightarrow K\text{-finite}$ and that in fact $\mathbf{Z}_I \vdash \text{strictly finite} \leftrightarrow \text{finite} \ \& \ \text{decidable}$. The strictly finite subsets of E correspond precisely to those which are bijective with initial segments of \mathbb{N} .

Frege's construction of the natural numbers can be carried out in \mathbf{Z}_I without the axiom of infinity, and the result shown to be equivalent to the postulation of the existence of a model of Peano's axioms, that is, the axiom of infinity. So we are led to define a *Frege structure* to be a pair (E, ν) with E a set and ν a function to E with domain a strictly inductive family F of subsets of E such that

$$\forall XY \in F. \nu(X) = \nu(Y) \leftrightarrow X \approx Y.^{17}$$

It can be shown that, for any Frege structure (E, ν) there is a subset N of E which is a model of Peano's axioms. To be precise, for $X \in \text{dom}(\nu) = F$ write $X^+ = X \cup \{\nu(X)\}$ and call a subfamily E of F *closed* if $\emptyset \in E$ and $X^+ \in E$ whenever $X \in E$ and $\nu(X) \notin X$. Let N be the intersection of all closed families, and define

$$\underline{0} = \nu(\emptyset), \quad N = \{\nu(X) : X \in N\}$$

and $s: N \rightarrow N$ by $s(\nu(X)) = \nu(X^+)$. Then $(N, s, \underline{0})$ is a model of Peano's axioms.

Conversely, each model $(N, s, 0)$ of Peano's axioms determines a Frege structure (N, ν) in which $\text{dom}(\nu)$ coincides with the family of (strictly) finite subsets of N .¹⁸ Here ν is given by

¹⁷ Here $X \approx Y$ stands for "there is a bijection between X and Y ".

¹⁸ Since \mathbb{N} is decidable, strict finiteness and finiteness of subsets of \mathbb{N} coincide.

$$v = \{(X, n) \in PN \times N: X \approx \{m:m < n\}\};$$

v assigns to each finite subset of \mathbb{N} the number of its elements.

Remark. In Frege's original formulation v was essentially a function from PE to E . Call such a Frege structure *full*. In classical set theory the natural number system determines a full Frege structure by defining, for $X \subseteq \mathbb{N}$,

$$v(X) = \begin{cases} |X| + 1 & \text{if } X \text{ is finite} \\ 0 & \text{if } X \text{ is infinite.} \end{cases}$$

But this cannot be the case in \mathbf{Z}_I , in view of the fact that *for any full Frege structure* (E, v) , *there is an injection* $\Omega \rightarrow E$. To see this, write $\underline{0} = v(\emptyset)$. Then for each $X, Y \subseteq \{\underline{0}\}$ we have

$$v(X) = v(Y) \leftrightarrow X \approx Y \leftrightarrow X = Y.$$

Thus the restriction of v to $P(\{\underline{0}\})$ is an injection into E , and since Ω is naturally isomorphic to $P(\{\underline{0}\})$, this determines an injection of Ω into E .

Therefore, if E is decidable, in particular if E is \mathbb{N} , Ω is also decidable, and **LEM** follows once again.

Classically, Zermelo-Fraenkel set theory **ZF** is obtained by adding to Zermelo set theory **Z** the axioms of *foundation* and *replacement*. Now the axiom of foundation asserts that each nonempty set u has a member x which is \in -*minimal*, that is, for which $x \cap u = \emptyset$. And it is easy to see that this implies **LEM**: an \in -minimal element of the set $\{0 \mid \alpha\} \cup \{1\}$ is either 0 or 1; if it is 0, α must hold, and if it is 1, α must fail; thus if foundation held we would get $\alpha \vee \neg\alpha$.

The appropriate substitute for the axiom of foundation is the scheme of \in -*induction*:

$$\mathbf{\in\text{-Ind}} \quad \forall x[\forall y \in x \alpha(y) \rightarrow \alpha(x)] \rightarrow \forall x \alpha(x).$$

Now *intuitionistic Zermelo-Fraenkel set theory* **ZF_I** is obtained by adding to the axioms of **Z_I** the scheme $\mathbf{\in\text{-Ind}}$ and the scheme of *replacement*

$$\mathbf{Rep} \quad \forall y \in x \exists ! z \alpha \rightarrow \exists w \forall y \in x \exists z \in w \alpha.$$

It is to be expected that the many classically equivalent definitions of *well-ordering* and *ordinal* become distinct with intuitionistic logic. The definitions we give here work reasonably well.

Definition. A set x is *transitive* if $\forall y \in x. y \subseteq x$; an *ordinal* is a transitive set of transitive sets. The class of ordinals is denoted by **Ord** and we use (italic) letters $\alpha, \beta, \gamma, \dots$ as variables ranging over it. A transitive subset of an ordinal is called a *subordinal*. An ordinal α is *simple* if $\forall \beta \gamma \in \alpha (\beta \in \gamma \vee \beta = \gamma \vee \gamma \in \beta)$.

Thus, for example, the ordinals 1, 2, 3, ... as well as the first infinite ordinal ω to be defined below, are all simple. Every subordinal of (hence every element) of a simple ordinal is simple. But, in contrast with classical set theory, intuitionistically not every ordinal can be simple, because the simplicity of the ordinal $\{0, \{0 \mid \alpha\}\}$ implies $\alpha \vee \neg\alpha$.

We next state the central properties of **Ord**.

Definition. The *successor* α^+ of an ordinal α is $\alpha \cup \{\alpha\}$; the *supremum* of a set A of ordinals is $\bigcup A$. The usual *order relations* are introduced on **Ord**:

$$\alpha < \beta \leftrightarrow \alpha \in \beta \quad \alpha \leq \beta \leftrightarrow \alpha \subseteq \beta.$$

It is now easily shown that successors and suprema of ordinals are again ordinals and that

$$\alpha < \beta \leftrightarrow \alpha^+ \leq \beta \quad \bigcup A \leq \beta \leftrightarrow \forall \alpha \in A. \alpha < \beta \leq \gamma \rightarrow \alpha < \gamma.$$

But straightforward arguments show that any of the following assertions (for arbitrary ordinals α, β, γ) implies LEM: (i) $\alpha < \beta \vee \alpha = \beta \vee \beta < \alpha$, (ii) $\alpha \leq \beta \vee \beta \leq \alpha$, (iii) $\alpha \leq \beta \rightarrow \alpha < \beta \vee \alpha = \beta$, (iv) $\alpha < \beta \rightarrow \alpha^+ < \beta \vee \alpha^+ = \beta$, (v) $\alpha \leq \beta < \gamma \rightarrow \alpha < \gamma$.

Definition. An ordinal α is a *successor* if $\exists \beta. \alpha = \beta^+$, a *weak limit* if $\forall \beta \in \alpha \exists \gamma \in \alpha. \beta \in \gamma$, and a *strong limit* if $\forall \beta \in \alpha. \beta^+ \in \alpha$.

Note that both the following assertions imply **LEM**: (i) every ordinal is zero, a successor, or a weak limit, (ii) all weak limits are strong limits. Assertion (i) follows from the observation that, for any formula α , if the specified disjunction applies to the ordinal $\{0 \mid \alpha\}$, then $\alpha \vee \neg \alpha$. As for assertion (ii), define

$$1_\alpha = \{0 \mid \alpha\}, 2_\alpha = \{0, 1_\alpha\}, \beta = \{0, 1_\alpha, 2_\alpha, 2_\alpha^+, 2_\alpha^{++}, \dots\}.$$

Then β is a weak limit, but a strong one only if $\alpha \vee \neg \alpha$.

As in classical set theory, in **ZF_I** a connection can be established between the class of ordinals and certain natural notions of well-founded or well-ordered structure. Thus a *well-founded* relation on a set A is a binary relation which is *inductive*, that is,

$$\forall X \subseteq A [\forall x \in A (\forall y < x. y \in X \rightarrow x \in X) \rightarrow A \subseteq X].$$

As for Foundation, the existence of *<-minimal* elements for any nontrivial relation $<$ implies **LEM**. But as in classical set theory, a well-founded relation has no infinite descending sequences and so is irreflexive. Moreover, the usual proof may be given in **ZF_I** to

justify *definitions by recursion* on a well-founded relation, so that we can make the following

Definition. If $<$ is a well-founded relation on a set A , the associated *rank function* $\rho_<: A \rightarrow \mathbf{Ord}$ is the (unique) function such that for each $x \in A$,

$$\rho_<(x) = \bigcup \{\rho_<(y)^+ : y < x\}.$$

When $<$ is \in restricted to an ordinal, it is easy to see that the associated rank function is the identity.

To obtain a characterization of the *order-types* represented by ordinals we make the following

Definition. A binary relation $<$ on a set A is *transitive* if $\forall xyz \in A (x < y \wedge y < z \rightarrow x < z)$, and *extensional* if $\forall xy \in A [\forall z (z < x \leftrightarrow z < y) \rightarrow x = y]$. A *well-ordering* is a transitive, extensional well-founded relation.

Now we can prove the

Theorem. The well-orderings are exactly those relations isomorphic to \in restricted to some ordinal.

Proof. It follows immediately from the axioms \in -**Ind** and **Ext** that the \in -relation well-orders every ordinal. Conversely, it is easy to prove by induction that the rank assigning function on any well-ordering is an isomorphism. ■

As observed above, we can justify definitions by \in -recursion on **Ord**, but we must avoid “taking cases” as is done classically.

Accordingly the definitions of *sums*, *products* and *exponentials* of ordinals have to be presented as *single equations*:

$$\begin{aligned}\alpha + \beta &= \alpha \cup \{\alpha + \delta : \delta \in \beta\} & \alpha \cdot \beta &= \{\alpha \cdot \delta + \gamma : \gamma \in \alpha, \delta \in \beta\} \\ \alpha^\beta &= 1 \cup \{\alpha^\delta \cdot \gamma + \varepsilon : \gamma \in \alpha, \delta \in \beta, \varepsilon \in \alpha^\delta\}.\end{aligned}$$

The *rank* $\text{rk}(x)$ of a set x is defined by recursion on \in by the equation $\text{rk}(x) = \bigcup\{\text{rk}(y)^+ : y \in x\}$. For $\alpha \in \mathbf{Ord}$ we define $V_\alpha = \bigcup\{P(V_\beta) : \beta < \alpha\}$. The rank function and the V_α have the following properties:

- (i) $\forall x \text{rk}(x) \in \mathbf{Ord}$
- (ii) $\forall \alpha \text{rk}(\alpha) = \alpha$
- (iii) $\forall x x \in V_{\text{rk}(x)+1}$
- (iv) $\alpha \leq \beta \rightarrow V_\alpha \subseteq V_\beta$
- (v) $x \subseteq y \in V_\alpha \rightarrow x \in V_\alpha$
- (vi) $V_\alpha \cap \mathbf{Ord} = \text{rk}(V_\alpha) \supseteq \alpha$.

All these are proved by routine induction arguments. In connection with (vi), we observe that the assertion $2 = V_2 \cap \mathbf{Ord}$ implies **LEM**. For by (v), $V_\alpha \cap \mathbf{Ord}$ is closed under subordinals, so in particular V_2 contains all the ordinals of the form $\{0 \mid \alpha\}$; but $\{0 \mid \alpha\} \in 2 \leftrightarrow \alpha \vee \neg\alpha$. In general $V_\alpha \cap \mathbf{Ord}$ can be very much bigger than α .

6. SMOOTH INFINITESIMAL ANALYSIS

Finally, we describe a remarkable new approach to infinitesimal analysis made possible by intuitionistic logic.

In the usual development of the calculus, for any differentiable function f on the real line \mathbf{R} , $y = f(x)$, it follows from Taylor's theorem that the increment $\delta y = f(x + \delta x) - f(x)$ in y attendant upon an increment δx in x is determined by an equation of the form

$$\delta y = f'(x)\delta x + A(\delta x)^2, \quad (1)$$

where $f'(x)$ is the derivative of $f(x)$ and A is a quantity whose value depends on both x and δx . Now if it were possible to take δx so *small* (but not demonstrably identical with 0) that $(\delta x)^2 = 0$ then (1) would assume the simple form

$$f(x + \delta x) - f(x) = \delta y = f'(x) \delta x. \quad (2)$$

We shall call a quantity having the property that its square is zero a *nilsquare infinitesimal* or simply an *infinitesimal* (or a *microquantity*). In *smooth infinitesimal analysis (SIA)* “enough” infinitesimals are present to ensure that equation (2) holds *nontrivially* for *arbitrary* functions $f: \mathbf{R} \rightarrow \mathbf{R}$. (Of course (2) holds trivially in standard mathematical analysis because there 0 is the sole infinitesimal in this sense.) The meaning of the term “nontrivial” here may be explicated in following way. If we replace

δx by the letter ε standing for an arbitrary infinitesimal, (2) assumes the form

$$f(x + \varepsilon) - f(x) = \varepsilon f'(x). \quad (3)$$

Ideally, we want the validity of this equation to be independent of ε , that is, given x , for it to hold for *all* infinitesimal ε . In that case the derivative $f'(x)$ may be *defined* as the unique quantity D such that the equation

$$f(x + \varepsilon) - f(x) = \varepsilon D$$

holds for all infinitesimal ε .

Setting $x = 0$ in this equation, we get in particular

$$f(\varepsilon) = f(0) + \varepsilon D, \quad (4)$$

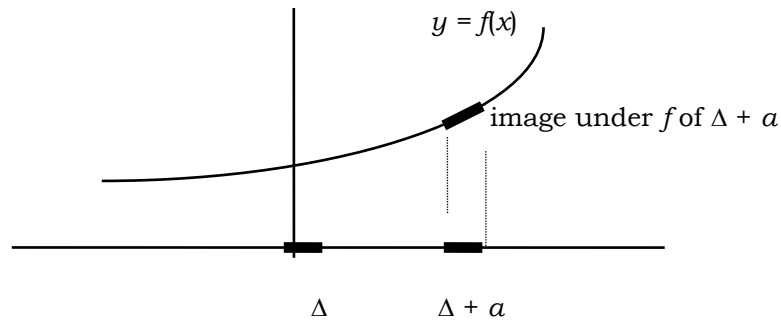
for all ε . *It is equation (4) that is taken as axiomatic in smooth infinitesimal analysis.* Let us write Δ for the set of infinitesimals, that is,

$$\Delta = \{x: x \in \mathbf{R} \wedge x^5 = 0\}.$$

Then it is postulated that, for any $f: \Delta \rightarrow \mathbf{R}$, there is a *unique* $D \in \mathbf{R}$ such that equation (4) holds for all ε . This says that the graph of f is a straight line passing through $(0, f(0))$ with slope D . Thus any function on Δ is what mathematicians term *affine*, and so this postulate is naturally termed the *principle of infinitesimal affinity*, or of *microstraightness*. It means that Δ *cannot be bent*

or broken: it is subject only to *translations and rotations*—and yet is not (as it would have to be in ordinary analysis) identical with a point. Δ may be thought of as an entity possessing position and attitude, but lacking true extension.

If we think of a function $y = f(x)$ as defining a curve, then, for any a , the image under f of the “infinitesimal interval” $\Delta + a$ obtained by translating Δ to a is straight and coincides with the tangent to the curve at $x = a$ (see figure immediately below). In this sense each curve is “infinitesimally straight”



From the principle of infinitesimal affineness we deduce the important

Principle of infinitesimal cancellation. *If $\varepsilon a = \varepsilon b$ for all ε , then $a = b$.*

For the premise asserts that the graph of the function $g: \Delta \rightarrow \mathbf{R}$ defined by $g(\varepsilon) = a\varepsilon$ has both slope a and slope b : the uniqueness condition in the principle of infinitesimal affineness then gives $a = b$. The principle of infinitesimal cancellation supplies the exact

sense in which there are “enough” infinitesimals in smooth infinitesimal analysis.

From the principle of infinitesimal cancellation it follows that Δ is *nondegenerate*, i.e. not identical with $\{0\}$. For if $\Delta = \{0\}$, we would have $\varepsilon.0 = \varepsilon.1$ for all ε , and infinitesimal cancellation would give $0 = 1$.

From the principle of infinitesimal affineness it also follows that *all functions on \mathbf{R} are continuous*, that is, *send neighbouring points to neighbouring points*. Here two points x, y on \mathbf{R} are said to be neighbours if $x - y$ is in Δ , that is, if x and y differ by an infinitesimal. To see this, given $f: \mathbf{R} \rightarrow \mathbf{R}$ and neighbouring points x, y , note that $y = x + \varepsilon$ with ε in Δ , so that

$$f(y) - f(x) = f(x + \varepsilon) - f(x) = \varepsilon f'(x).$$

But clearly any multiple of an infinitesimal is also an infinitesimal, so $\varepsilon f'(x)$ is infinitesimal, and the result follows.

In fact, since equation (3) holds for any f , it also holds for its derivative f' ; it follows that functions in smooth infinitesimal analysis are differentiable arbitrarily many times, thereby justifying the use of the term “smooth”.

Let us derive a basic law of the differential calculus, the *product rule*:

$$(fg)' = f'g + fg'.$$

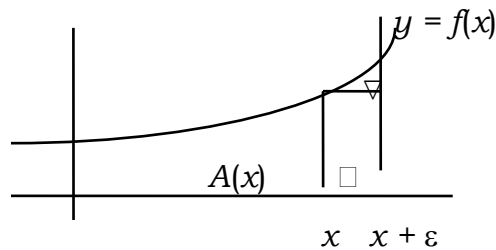
To do this we compute

$$(fg)(x + \varepsilon) = (fg)(x) + (fg)'(x) = f(x)g(x) + (fg)'(x),$$

$$\begin{aligned} (fg)(x + \varepsilon) &= f(x + \varepsilon)g(x + \varepsilon) = [f(x) + f'(x)\varepsilon][g(x) + g'(x)\varepsilon] \\ &= f(x)g(x) + \varepsilon(f'g + fg') + \varepsilon^2 f'g' \\ &= f(x)g(x) + \varepsilon(f'g + fg'), \end{aligned}$$

since $\varepsilon^2 = 0$. Therefore $\varepsilon(fg)' = \varepsilon(f'g + fg')$, and the result follows by infinitesimal cancellation.

Next, we derive the *Fundamental Theorem of the Calculus*.



Let J be a closed interval $[a, b] = \{x: a \leq x \leq b\}$ in \mathbf{R} and $f: J \rightarrow \mathbf{R}$; let $A(x)$ be the area under the curve $y = f(x)$ as indicated above. Then, using equation (3),

$$\varepsilon A'(x) = A(x + \varepsilon) - A(x) = \square + \nabla = \varepsilon f(x) + \nabla.$$

Now by infinitesimal affineness ∇ is a triangle of area $\frac{1}{2} \varepsilon \cdot \varepsilon f'(x) = 0$.

Hence $\varepsilon A'(x) = \varepsilon f(x)$, so that, by infinitesimal cancellation,

$$A'(x) = f(x).$$

We observe that the postulates of smooth infinitesimal analysis are *incompatible with the law of excluded middle of classical logic*. This incompatibility can be demonstrated in two ways, one informal and the other rigorous. First the informal argument. Consider the function f defined for real numbers x by $f(x) = 1$ if $x = 0$ and $f(x) = 0$ whenever $x \neq 0$. If the law of excluded middle held, each real number would then be either equal or unequal to 0, so that the function f would be defined on the whole of \mathbf{R} . But, considered as a function with domain \mathbf{R} , f is clearly discontinuous. Since, as we know, in smooth infinitesimal analysis every function on \mathbf{R} is continuous, f cannot have domain \mathbf{R} there¹⁹. So the law of excluded middle fails in smooth infinitesimal analysis. To put it succinctly, *universal continuity implies the failure of the law of excluded middle*.

Here now is the rigorous argument. We show that the failure of the law of excluded middle can be derived from the principle of infinitesimal cancellation. To begin with, if $x \neq 0$, then $x^2 \neq 0$, so that, if $x^2 = 0$, then necessarily not $x \neq 0$. This means that

$$\text{for all infinitesimal } \varepsilon, \text{ not } \varepsilon \neq 0.$$

(*)

Now suppose that the law of excluded middle were to hold. Then we would have, for any ε , either $\varepsilon = 0$ or $\varepsilon \neq 0$. But (*) allows us to

¹⁹ The domain of f is in fact $(\mathbf{R} - \{0\}) \cup \{0\}$, which, because of the failure of the law of excluded middle in **SIA**, is provably unequal to \mathbf{R} .

eliminate the second alternative, and we infer that, for all ε , $\varepsilon = 0$. This may be written

$$\text{for all } \varepsilon, \varepsilon.1 = \varepsilon.0,$$

from which we derive by infinitesimal cancellation the falsehood $1 = 0$. So again the law of excluded middle must fail.

The “internal” logic of smooth infinitesimal analysis is accordingly not full classical logic. It is, instead, *intuitionistic* logic. In our brief sketch we did not notice this “change of logic” because, like much of elementary mathematics, the topics we discussed are naturally treated by constructive means such as direct computation.

ALGEBRAIC AND ORDER STRUCTURE OF \mathbf{R}

What are the *algebraic* and *order structures* on \mathbf{R} in **SIA**? As far as the former is concerned, there is little difference from the classical situation: in **SIA** \mathbf{R} is equipped with the usual addition and multiplication operations under which it is a field. In particular, \mathbf{R} satisfies the condition that each $x \neq 0$ has a multiplicative inverse. Notice, however, that since in **SIA** no microquantity (apart from 0 itself) is provably $\neq 0$, microquantities are not required to have multiplicative inverses (a requirement which would lead to inconsistency). From a strictly algebraic standpoint, \mathbf{R} in **SIA** differs from its classical counterpart only in being required to satisfy the principle of infinitesimal cancellation.

The situation is different, however, as regards the order structure of \mathbf{R} in **SIA**. Because of the failure of the law of excluded

middle, the order relation $<$ on \mathbf{R} in SIA cannot satisfy the trichotomy law

$$x < y \vee y < x \vee x = y,$$

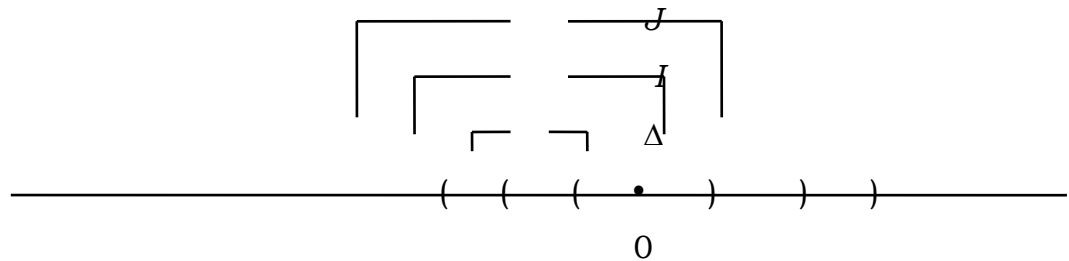
and accordingly $<$ must be a *partial*, rather than a *total* ordering. Since microquantities do not have multiplicative inverses, and \mathbf{R} is a field, any microquantity ε must satisfy

$$\neg \varepsilon < 0 \wedge \neg \varepsilon > 0.$$

Accordingly, if we define the relation \leq (“not less than”) $x < y$, then, for any microquantity ε we have

$$\varepsilon \leq 0 \wedge \varepsilon \geq 0.$$

Using these ideas we can identify three distinct *infinitesimal neighbourhoods* of 0 on \mathbf{R} in SIA, each of which is included in its successor.



First, the set Δ of microquantities itself, next, the set $I = \{x \in \mathbf{R}: \neg x \neq 0\}$ of elements indistinguishable from 0; finally, the set $J = \{x \in \mathbf{R}: x \leq 0 \wedge x \geq 0\}$ of elements neither less nor greater than 0.

These three may be thought of as the infinitesimal neighbourhoods of 0 defined *algebraically*, *logically*, and *order-theoretically*, respectively. Observe that none of these is degenerate.

SIA VERSUS CONSTRUCTIVE ANALYSIS

SIA may be furnished with the following *axiomatic description*.

Axioms for the continuum, or smooth real line \mathbf{R} . These are the usual axioms for a field expressed in terms of two operations $+$ and \cdot and two distinguished elements 0, 1. In particular every nonzero element of \mathbf{R} is invertible.

Axioms for the strict order relation $<$ on \mathbf{R} . These are:

1. $a < b$ and $b < c$ implies $a < c$.
2. $\neg(a < a)$
3. $a < b$ implies $a + c < b + c$ for any c .
4. $a < b$ and $0 < c$ implies $a \cdot c < b \cdot c$.
5. either $0 < a$ or $a < 1$.

The subset $\Delta = \{x: x^2 = 0\}$ of \mathbf{R} is subject to the

Infinitesimal Affineness Principle. *For any map $g: \Delta \rightarrow \mathbf{R}$ there exist unique $a, b \in \mathbf{R}$ such that, for all ε , we have*

$$g(\varepsilon) = a + b \cdot \varepsilon.$$

From these three axioms it follows that the continuum in **SIA** differs in certain key respects from its counterpart in *constructive analysis* **CA**, which was introduced in Chapter 1. To begin with, a basic property of the strict ordering relation $<$ in **CA**, namely,

$$(*) \quad \neg(x < y \vee y < x) \rightarrow x = y$$

is incompatible with the axioms of **SIA**. For (*) implies

$$(**) \quad \forall x \neg(x < 0 \vee 0 < x) \rightarrow x = 0.$$

Thus in **CA** the set Δ of infinitesimals would be degenerate (i.e., identical with $\{0\}$), while, as we have seen, the nondegeneracy of Δ in **SIA** is one of its characteristic features.

Next, call a binary relation S on **R** *stable* if it satisfies

$$\forall x \forall y (\neg\neg x R y \rightarrow x R y).$$

As we have observed, in **CA**, the equality relation is stable. But in **SIA** it is not stable, for, if it were, I would be degenerate, which we have observed is not the case in **SIA**.

INDECOMPOSABILITY OF THE CONTINUUM IN **SIA**

A *stationary point* a in **R** of a function $f: \mathbf{R} \rightarrow \mathbf{R}$ is defined to be one in whose vicinity “infinitesimal variations” fail to change the value of f , that is, such that $f(a + \varepsilon) = f(a)$ for all ε . This means that

$f(a) + \varepsilon f'(a) = f(a)$, so that $\varepsilon f'(a) = 0$ for all ε , whence it follows from infinitesimal cancellation that $f'(a) = 0$. This is *Fermat's rule*.

An important postulate concerning stationary points that we adopt in smooth infinitesimal analysis is the

Constancy Principle. If every point in an interval J is a stationary point of $f: J \rightarrow \mathbf{R}$ (that is, if f' is identically 0), then f is constant.

Put succinctly, “universal local constancy implies global constancy”. It follows from this that two functions with identical derivatives differ by at most a constant.

In ordinary analysis the continuum \mathbf{R} and all closed intervals are connected in the sense that they cannot be split into two non empty subsets neither of which contains a limit point of the other. In smooth infinitesimal analysis they have the vastly stronger property of *indecomposability*: they cannot be split *in any way whatsoever* into two disjoint nonempty subsets. For suppose $\mathbf{R} = U \cup V$ with $U \cap V = \emptyset$. Define $f: \mathbf{R} \rightarrow \{0, 1\}$ by $f(x) = 1$ if $x \in U$, $f(x) = 0$ if $x \in V$. We claim that f is constant. For we have

$$(f(x) = 0 \text{ or } f(x) = 1) \quad \& \quad (f(x + \varepsilon) = 0 \text{ or } f(x + \varepsilon) = 1).$$

This gives 4 possibilities:

- (i) $f(x) = 0 \quad \& \quad f(x + \varepsilon) = 0$
- (ii) $f(x) = 0 \quad \& \quad f(x + \varepsilon) = 1$
- (iii) $f(x) = 1 \quad \& \quad f(x + \varepsilon) = 0$
- (iv) $f(x) = 1 \quad \& \quad f(x + \varepsilon) = 1$

Possibilities (ii) and (iii) may be ruled out because f is continuous. This leaves (i) and (iv), in either of which $f(x) = f(x + \varepsilon)$. So f is locally, and hence globally, constant, that is, constantly 1 or 0. In the first case $V = \emptyset$, and in the second $U = \emptyset$. The argument for an arbitrary closed interval is similar.

From the indecomposability of closed intervals it follows that *all intervals in \mathbf{R} are indecomposable*. To do this we employ the following

Lemma. Suppose that A is an inhabited²⁰ subset of \mathbf{R} satisfying

(*) for any $x, y \in A$ there is an indecomposable set B such that

$$\{x, y\} \subseteq B \subseteq A.$$

Then A is indecomposable.

Proof. Suppose A satisfies (*) and $A = U \cup V$ with $U \cap V = \emptyset$. Since A is inhabited, we may choose $a \in A$. Then $a \in U$ or $a \in V$. Suppose $a \in U$; then if $y \in V$ there is an indecomposable B for which $\{a, y\} \subseteq B \subseteq A = U \cup V$. It follows that $B = (B \cap U) \cup (B \cap V)$, whence $B \cap U = \emptyset$ or $B \cap V = \emptyset$. The former possibility is ruled out by the fact that $a \in B \cap U$, so $B \cap V = \emptyset$, contradicting $y \in B \cap V$. Therefore $y \in V$ is impossible; since this is the case for arbitrary y , we conclude that $V = \emptyset$. Similarly, if $a \in V$, then $U = \emptyset$, so that A is indecomposable as claimed.

²⁰ A set A is *inhabited* if $\exists x. x \in A$.

We use this lemma to show that the open interval $(0, 1) = \{x \in \mathbf{R} : 0 < x < 1\}$ is indecomposable; similar arguments work for arbitrary intervals. In fact, if $\{x, y\} \subseteq (0, 1)$, it is easy to verify that

$$\{x, y\} \subseteq [xy/x+y, 1-xy/2-x-y] \subseteq (0, 1).$$

Thus, in view of the indecomposability of closed intervals, $(0, 1)$ satisfies condition (*) of the lemma, and so is indecomposable.

In some versions of **SIA** the ordering of \mathbf{R} is subject to the *axiom of distinguishability*:

$$(*) \quad x \neq y \rightarrow x < y \vee y < x.$$

Aside from certain infinitesimal subsets to be discussed below, in these versions of **SIA** indecomposable subsets of \mathbf{R} correspond to connected subsets of \mathbf{R} in classical analysis, that is, to intervals. In particular, in versions of **SIA** subject to (*) any puncturing of \mathbf{R} is *decomposable*, for it follows immediately from (*) that

$$\mathbf{R} - \{a\} = \{x : x > a\} \cup \{x : x < a\}.$$

Similarly, the set $\mathbf{R} - \mathbf{Q}$ of irrational numbers is decomposable as

$$\mathbf{R} - \mathbf{Q} = [\{x : x > 0\} - \mathbf{Q}] \cup [\{x : x < 0\} - \mathbf{Q}].$$

This is in sharp contrast with the situation in *intuitionistic analysis* **IA**, that is, **CA** augmented by certain principles (Kripke's scheme, the continuity principle, and bar induction). For in **IA** not only is any puncturing of \mathbf{R} indecomposable, but that this is even the case for the set of irrational numbers. This would seem to indicate that

in some sense the continuum in **SIA** is considerably less “syrupy”²¹ than its counterpart in **IA**.

It can be shown that the various infinitesimal neighbourhoods of 0 are indecomposable. For example, the indecomposability of Δ can be established as follows. Suppose $f: \Delta \rightarrow \{0, 1\}$. Then by Microaffineness there are unique $a, b \in \mathbf{R}$ such that $f(\varepsilon) = a + b.\varepsilon$ for all ε . Now $a = f(0) = 0$ or 1; if $a = 0$, then $b.\varepsilon = f(\varepsilon) = 0$ or 1, and clearly $b.\varepsilon \neq 1$. So in this case $f(\varepsilon) = 0$ for all ε . If on the other hand $a = 1$, then $1 + b.\varepsilon = f(\varepsilon) = 0$ or 1; but $1 + b.\varepsilon = 0$ would imply $b.\varepsilon = -1$ which is again impossible. So in this case $f(\varepsilon) = 1$ for all ε . Therefore f is constant and Δ indecomposable.

In **SIA** *nilpotent infinitesimals* are defined to be the members of the sets

$$\Delta_k = \{x \in \mathbf{R}: x^{k+1} = 0\},$$

for $k = 1, 2, \dots$, each of which may be considered an infinitesimal neighbourhood of 0. These are subject to the

Micropolynomiality Principle. *For any $k \geq 1$ and any $g: \Delta_k \rightarrow \mathbf{R}$, there exist unique $a, b_1, \dots, b_k \in \mathbf{R}$ such that for all $\delta \in \Delta_k$ we have*

$$g(\delta) = a + b_1\delta + b_2\delta^2 + \dots + b_k\delta^k.$$

Micropolynomiality implies that no Δ_k coincides with $\{0\}$.

An argument similar to that establishing the indecomposability of Δ does the same for each Δ_k . Thus let $f: \Delta_k \rightarrow \{0, 1\}$; Micropolynomiality implies the existence of $a, b_1, \dots, b_k \in \mathbf{R}$ such that $f(\delta) = a + \zeta(\delta)$, where $\zeta(\delta) = b_1\delta + b_2\delta^2 + \dots + b_k\delta^k$. Notice that $\zeta(\delta) \in \Delta_k$, that is, $\zeta(\delta)$ is nilpotent. Now $a = f(0) = 0$ or 1; if $a = 0$ then $\zeta(\delta) = f(\delta) = 0$ or 1, but since $\zeta(\delta)$ is nilpotent it cannot = 1. Accordingly in this case $f(\delta) = 0$ for all $\delta \in \Delta_k$. If on the other hand a

²¹ It should be emphasized that this phenomenon is a consequence of (*): it cannot necessarily be affirmed in versions of **SIA** not including this axiom.

$= 1$, then $1 + \zeta(\delta) = f(\delta) = 0$ or 1 , but $1 + \zeta(\delta) = 0$ would imply $\zeta(\delta) = -1$ which is again impossible. Accordingly f is constant and Δ_k indecomposable.

The union \mathbf{D} of all the Δ_k is the *set of nilpotent infinitesimals*, another infinitesimal neighbourhood of 0 . The indecomposability of \mathbf{D} follows immediately by applying the Lemma above.

The next infinitesimal neighbourhood of 0 is the closed interval $[0, 0]$, which, as a closed interval, is indecomposable. It is easily shown that $[0, 0]$ includes \mathbf{D} , so that it does not coincide with $\{0\}$.

It is also easily shown, using axioms 2 and 6, that $[0, 0]$ coincides with the set

$$\mathbf{I} = \{x \in \mathbf{R} : \neg\neg x = 0\}.$$

So \mathbf{I} is indecomposable. (In fact the indecomposability of \mathbf{I} can be proved independently of axioms 1-6 through the general observation that, if A is indecomposable, then so is the set $A^* = \{x : \neg\neg x \in A\}$.)

Finally, we observe that the sequence of infinitesimal neighbourhoods of 0 generates a strictly ascending sequence of decomposable subsets containing $\mathbf{R} - \{0\}$, namely:

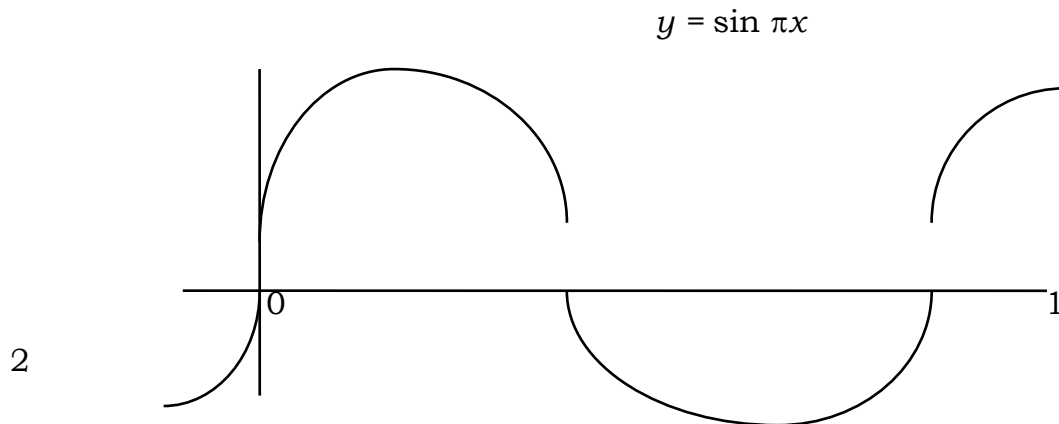
$$\begin{aligned} \mathbf{R} - \{0\} \subset (\mathbf{R} - \{0\}) \cup \{0\} \subset (\mathbf{R} - \{0\}) \cup \Delta_1 \subset (\mathbf{R} - \{0\}) \cup \Delta_2 \subset \dots (\mathbf{R} - \{0\}) \cup \\ \mathbf{D} \subset \\ (\mathbf{R} - \{0\}) \cup [0, 0]. \end{aligned}$$

NATURAL NUMBERS AND INVERTIBLE INFINITESIMALS IN **SIA**

In certain models of **SIA** the system of *natural numbers* possesses some subtle and intriguing features which make it

possible to introduce another type of infinitesimal—the so-called *invertible* infinitesimals—resembling those of nonstandard analysis, whose presence engenders yet another infinitesimal neighbourhood of 0 properly containing all those introduced above.

In **SIA** the set **N** of natural numbers can be defined to be the smallest subset of **R** which contains 0 and is closed under the operation of adding 1. In some models of **SIA**, **R** satisfies the *Archimedean principle* that every real number is majorized by a natural number. However, models of **SIA** have been constructed in which **R** is not Archimedean in this sense. In these models it is more natural to consider, in place of **N**, the set **N*** of *smooth natural numbers*



defined by

$$\mathbf{N}^* = \{x \in \mathbf{R}: 0 \leq x \wedge \sin \pi x = 0\}.$$

\mathbf{N}^* is the set of points of intersection of the smooth curve $y = \sin \pi x$ with the positive x -axis. In these models \mathbf{R} can be shown to possess the Archimedean property *provided that in the definition \mathbf{N} is replaced by \mathbf{N}^** . In these models, then, \mathbf{N} is a proper subset of \mathbf{N}^* : the members of $\mathbf{N}^* - \mathbf{N}$ may be considered *nonstandard integers*. Multiplicative inverses of nonstandard integers are infinitesimals, but, being themselves invertible, they are of a different type from the ones we have considered so far. It is quite easy to show that they, as well as the infinitesimals in J (and so also those in Δ and I) are all contained in the set—a further infinitesimal neighbourhood of 0—

$$K = \{x \in \mathbf{R}: \forall n \in \mathbf{N}. -1/n+1 < x < 1/n+1\}$$

of *infinitely small* elements of \mathbf{R} . The members of the set

$$In = \{x \in K: x \neq 0\}$$

of invertible elements of K are naturally identified as *invertible* infinitesimals. Being obtained as inverses of “infinitely large” reals (i.e. reals r satisfying $\forall n \in \mathbf{N}. n < r \vee \forall n \in \mathbf{N}. r < -n$) the members of In are the counterparts in **SIA** of the infinitesimals of nonstandard analysis.

In the past physicists showed no hesitation in employing infinitesimal methods²², the use of which in turn relied on the implicit assumption that the (physical) world is smooth, or at least that the maps encountered there are differentiable as many times as needed. For this reason smooth infinitesimal analysis (SIA) provides an ideal framework for the rigorous derivation of results in classical physics.

The notions and principles of SIA we will use in our derivations are the following:

- The domain \mathbf{R} of reals contains the nondegenerate set $\Delta = \{x: x^2 = 0\}$ of *microquantities*. We use letters $\varepsilon, \eta, \zeta, \tau$ as variables ranging over Δ .
- Δ is subject to the *principle of microaffineness*: any map $f: \Delta \rightarrow \mathbf{R}$ is affine, that is, for some constant a ,

$$f(\varepsilon) = f(0) + a\varepsilon.$$

As a consequence the image of Δ under any map is straight.

- For any map $f: \mathbf{R} \rightarrow \mathbf{R}$, the derivative f' of f is related to f by the identity in x and ε

$$f(x + \varepsilon) = f(x) + \varepsilon f'(x).$$

²² In this connection we recall the words of Hermann Weyl:

The principle of gaining knowledge of the external world from the behaviour of its infinitesimal parts is the mainspring of the theory of knowledge in infinitesimal physics as in Riemann's geometry, and, indeed, the mainspring of all the eminent work of Riemann (1922, p. 92).

If (as Hilbert said) set theory is "Cantor's paradise" then smooth infinitesimal analysis is nothing less than "Riemann's paradise".

- Given a function $f: \mathbf{R}^n \rightarrow \mathbf{R}$ of n variables x_1, \dots, x_n , the partial derivative $\frac{\partial f}{\partial x_i}$ is defined as usual to be the derivative of the function $f(a_1, \dots, x_i, \dots, a_n)$ obtained by fixing the values of all the variables apart from x_i . In that case, for an arbitrary microquantity ε , we have

$$f(x_1, \dots, x_i + \varepsilon, \dots, x_n) = f(x_1, \dots, x_n) + \varepsilon \frac{\partial f}{\partial x_i}(x_1, \dots, x_n).$$

Using the fact that $\varepsilon^2 = 0$, it is then easily shown that

$$f(x_1 + a_1\varepsilon, \dots, x_n + a_n\varepsilon) = f(x_1, \dots, x_n) + \varepsilon \sum_{i=1}^n a_i \frac{\partial f}{\partial x_i}(x_1, \dots, x_n).$$

- We also have the *Principle of Microcancellation*, viz..

$$\text{for } a, b \in \mathbf{R}, \forall \varepsilon [\varepsilon a = \varepsilon b] \Rightarrow a = b.$$

Microcancellation is a rigorous version of the process, familiar to physicists and engineers, of cancellation of differentials, as in

$$adx = bdx \Rightarrow ax = \int adx = \int bdx = bx \Rightarrow a = b.$$

The Heat Equation

$$\frac{\mathbf{S}}{O \quad P \quad Q}$$

Suppose we are given a heated wire W ; let $T(x,t)$ be the temperature at the point P at distance x along W from some given point O on it at time t .

Consider the segment S of W extending from P to the point Q at distance ε from P . The heat content of S is $k\varepsilon T_{\text{average}}$, where T_{average} is the average temperature over S and k is a constant depending on the material of the wire. Assuming that T_{average} is a convex combination $\lambda T(x+\varepsilon, t) + (1 - \lambda)T(x, t)$ of the temperatures at the endpoints of S , we see that the heat content of S is

$$\begin{aligned} k\varepsilon T_{\text{average}} &= k\varepsilon[\lambda T(x+\varepsilon, t) + (1 - \lambda)T(x, t)] \\ &= k\varepsilon[\lambda(T(x, t) + \varepsilon \frac{\partial T}{\partial x}) + (1 - \lambda)T(x, t)] \\ &= k\varepsilon T(x, t) \quad (\text{noting that } \varepsilon^2 = 0). \end{aligned}$$

Accordingly the change in heat content in S from time t to time $t + \eta$ is

$$(*) \quad k\varepsilon[T(x, t + \eta) - T(x, t)] = k\varepsilon\eta \frac{\partial T}{\partial t}(x, t).$$

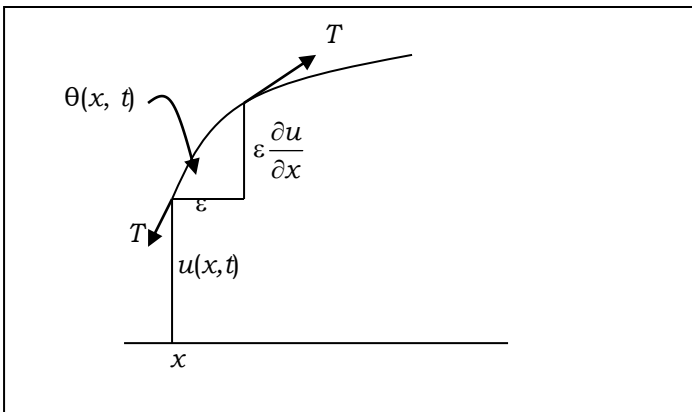
On the other hand, the rate of thermal flow across P is proportional to the temperature there, and so equal to $m \frac{\partial T}{\partial x}(x, t)$, where m is a constant depending on the material of the wire. Similarly, the rate of thermal flow across the point Q is $m \frac{\partial T}{\partial x}(x + \varepsilon, t)$. Thus the thermal transfer across P from time t to time $t + \eta$ is $m\eta \frac{\partial T}{\partial x}(x, t)$ and that across Q is $m\eta \frac{\partial T}{\partial x}(x + \varepsilon, t)$. So the net change in heat content in S from time t to time $t + \eta$ is

$$m\eta\left[\frac{\partial T}{\partial x}(x+\varepsilon, t) - \frac{\partial T}{\partial x}(x, t)\right] = m\eta\varepsilon\frac{\partial^2 T}{\partial x^2}(x, t).$$

Equating the rhs of this with the rhs of (*), cancelling η and ε , and setting $c = k/m$ yields the heat equation

$$\frac{\partial^2 T}{\partial x^2} = c \frac{\partial T}{\partial t}.$$

The Wave Equation



Assume that the tension T and density ρ of a stretched string are both constant throughout its length (and independent of the time). Let $u(x, t)$, $\theta(x, t)$ be, respectively, the vertical displacement of the string and the angle between the string and the horizontal at position x and time t .

Consider a microelement of the string between x and $x + \varepsilon$ at time t . Its mass is $\varepsilon\rho\cos\theta(x, t)$ and its vertical acceleration $\frac{\partial^2 u}{\partial t^2}(x, t)$. The vertical force on the element is

$$T[\sin \theta(x + \varepsilon, t) - \sin \theta(x, t)] = \varepsilon T \cos \theta(x, t) \frac{\partial \theta}{\partial x}(x, t).$$

By Newton's second law, we may equate the force with mass \times acceleration giving

$$\varepsilon\rho\cos\theta\frac{\partial^2 u}{\partial t^2} = \varepsilon T\cos\theta\frac{\partial\theta}{\partial x}.$$

Cancelling the universally quantified ε gives

$$\rho\cos\theta\frac{\partial^2 u}{\partial t^2} = T\cos\theta\frac{\partial\theta}{\partial x}.$$

Since $\cos\theta \neq 0$ it may also be cancelled to give

$$(1) \quad \rho\frac{\partial^2 u}{\partial t^2} = T\frac{\partial\theta}{\partial x}.$$

Now we recall the fundamental equation governing sines and cosines

.

$$(2) \quad \sin\theta = \cos\theta \frac{\partial u}{\partial x}.$$

Applying $\frac{\partial}{\partial x}$ to both sides of this gives

$$\cos\theta \frac{\partial\theta}{\partial x} = -\sin\theta \frac{\partial\theta}{\partial x} \frac{\partial u}{\partial x} + \cos\theta \frac{\partial^2 u}{\partial x^2}.$$

Substituting (2) in this latter equation yields

$$\cos\theta \frac{\partial\theta}{\partial x} = -\cos\theta \frac{\partial\theta}{\partial x} \left(\frac{\partial u}{\partial x}\right)^2 + \cos\theta \frac{\partial^2 u}{\partial x^2}.$$

Cancelling $\cos\theta$ and rearranging gives

$$\frac{\partial\theta}{\partial x} = \frac{\partial^2 u}{\partial x^2} / 1 + \left(\frac{\partial u}{\partial x}\right)^2.$$

Substituting this in (1) yields the rigorous wave equation

$$(3) \quad \frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2} / 1 + \left(\frac{\partial u}{\partial x}\right)^2$$

with $c = \sqrt{\frac{T}{\rho}}$.

When the amplitude of vibration is small we may assume that

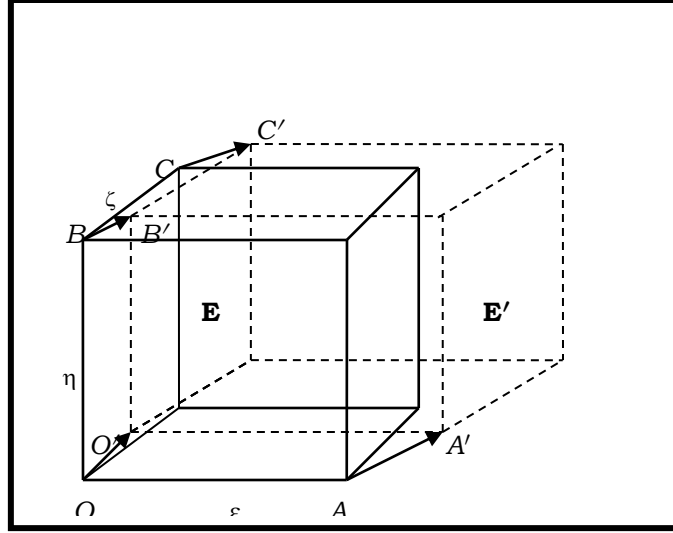
$\left(\frac{\partial u}{\partial x}\right)^2 = 0$ and in that case (3) becomes the familiar wave equation

$$\frac{\partial^2 u}{\partial t^2} = c^2 \frac{\partial^2 u}{\partial x^2}.$$

Euler's Equation of Continuity for Fluids.

In his derivation of the equation Euler employs infinitesimal volume elements of sufficient minuteness so as to preserve their rectilinear *shape* under infinitesimal flow, yet allowing their *volume* to undergo infinitesimal change. This idea was to become fundamental in continuum mechanics. The derivation in SIA will follow Euler's very closely, but the use of microquantities and the microcancellation axiom will make the argument entirely rigorous.

Here we are given a fluid free of viscosity but of varying density flowing smoothly in space. At any point $O = (x, y, z)$ in the fluid and at any time t , the fluid's density ρ and the components u , v , w of the fluid's velocity are given as functions of x , y , z , t . Following Euler, we consider the elementary volume element **E**—a microparallelepiped—with origin O and edges OA , OB , OC of microlengths ε , η , ζ and so of mass $\varepsilon\eta\zeta\rho$:



Fluid flow during the microtime τ transforms the volume element \mathbf{E} into the microparallelepiped \mathbf{E}' with vertices O' , A' , B' , C' . We first calculate the length of the side $O'A'$. Now the rate at which A is moving away from O in the x -direction is

$$u(x + \varepsilon, y, z, t) - u(x, y, z, t) = \varepsilon \frac{\partial u}{\partial x}.$$

The change in length of OA during the microtime τ is thus $\varepsilon \tau \frac{\partial u}{\partial x}$, so that the length of $O'A'$ is $\varepsilon + \varepsilon \tau \frac{\partial u}{\partial x} = \varepsilon \left(1 + \tau \frac{\partial u}{\partial x} \right)$. Similarly, the lengths of $O'B'$ and $O'C'$ are, respectively,

$$\eta \left(1 + \tau \frac{\partial v}{\partial y} \right), \quad \zeta \left(1 + \tau \frac{\partial w}{\partial z} \right).$$

The volume of \mathbf{E}' is the product of these three quantities, which, using the fact that $\tau^2 = 0$, comes out as

$$(A) \quad \varepsilon\eta\zeta \left[1 + \tau \left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} + \frac{\partial w}{\partial z} \right) \right].$$

Since the coordinates of O' are $(x+u\tau, y+v\tau, z+w\tau)$, the fluid density ρ' there at time $t + \tau$ is, using (2),

$$(B) \quad \rho + \tau \left(\frac{\partial \rho}{\partial t} + u \frac{\partial \rho}{\partial x} + v \frac{\partial \rho}{\partial y} + w \frac{\partial \rho}{\partial z} \right).$$

The mass of \mathbf{E}' is then the product of (A) and (B), which, again using the fact that $\tau^2 = 0$, comes out as

$$(C) \quad \varepsilon\eta\zeta\rho + \varepsilon\eta\zeta\tau \left(\frac{\partial \rho}{\partial t} + \rho \frac{\partial u}{\partial x} + \rho \frac{\partial v}{\partial y} + \rho \frac{\partial w}{\partial z} + u \frac{\partial \rho}{\partial x} + v \frac{\partial \rho}{\partial y} + w \frac{\partial \rho}{\partial z} \right).$$

Now by the principle of conservation of mass, the masses of the fluid in \mathbf{E} and \mathbf{E}' are the same, so equating the mass $\varepsilon\eta\zeta\rho$ of \mathbf{E} to the mass of \mathbf{E}' given by (C) yields

$$\varepsilon\eta\zeta\tau \left(\frac{\partial \rho}{\partial t} + \rho \frac{\partial u}{\partial x} + \rho \frac{\partial v}{\partial y} + \rho \frac{\partial w}{\partial z} + u \frac{\partial \rho}{\partial x} + v \frac{\partial \rho}{\partial y} + w \frac{\partial \rho}{\partial z} \right) = 0.$$

Microcancellation gives

$$\frac{\partial \rho}{\partial t} + \rho \frac{\partial u}{\partial x} + \rho \frac{\partial v}{\partial y} + \rho \frac{\partial w}{\partial z} + u \frac{\partial \rho}{\partial x} + v \frac{\partial \rho}{\partial y} + w \frac{\partial \rho}{\partial z} = 0,$$

i.e.,

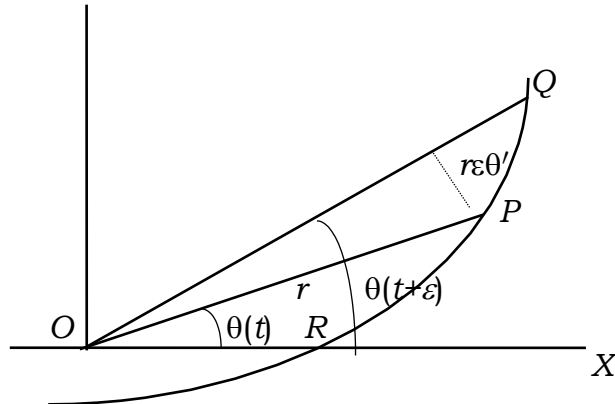
$$\frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x}(\rho u) + \frac{\partial}{\partial y}(\rho v) + \frac{\partial}{\partial z}(\rho w) = 0,$$

Euler's equation of continuity.

The Kepler-Newton areal law of motion under a central force.

We suppose that a particle executes plane motion under the influence of a force directed towards some fixed point O . If P is a point on the particle's trajectory with coordinates x, y , we write r for the length of the line PO and θ for the angle that it makes with the x -axis OX . Let A be the area of the sector ORP , where R is the point of intersection of the trajectory with OX . We regard x, y, r, θ as functions of a time variable t : thus

$$x = x(t), y = y(t), r = r(t), \theta = \theta(t), A = A(t).$$



Now let Q be a point on the trajectory at which the time variable has value $t + \varepsilon$, with ε in Δ . Then by Microaffineness the sector OPQ is a triangle of base $r(t + \varepsilon) = r + \varepsilon r'$ and height

$$r \sin[\theta(t + \varepsilon) - \theta(t)] = r \sin \varepsilon \theta' = r \varepsilon \theta'.^{23}$$

The area of OPQ is accordingly

$$2 \text{ base} \times \text{height} = 2 (r + \varepsilon r') r \varepsilon \theta' = 2(r^2 \varepsilon \theta' + \varepsilon^2 r r' \theta') = 2 r^2 \varepsilon \theta'.$$

Therefore

$$\varepsilon A'(t) = A(t + \varepsilon) - A(t) = \text{area } OPQ = 2\varepsilon r^2 \theta',$$

so that, cancelling ε ,

$$A'(t) = 2r^2 \theta'. \quad (*)$$

Now let $H = H(t)$ be the acceleration towards O induced by the force. Resolving the acceleration along and normal to OX , we have

$$x'' = H \cos \theta \quad y'' = H \sin \theta.$$

²³ Here we note that $\sin \varepsilon = \varepsilon$ for microquantities ε : recall that $\sin x$ is approximately equal to x for small values of x .

Also $x = r \cos\theta$, $y = r \sin\theta$. Hence

$$yx'' = Hy \cos\theta = Hr \sin\theta \cos\theta \quad xy'' = Hx \sin\theta = Hr \sin\theta \cos\theta,$$

from which we infer that

$$(xy' - yx')' = xy'' - yx'' = 0.$$

Hence

$$xy' - yx' = k, \tag{**}$$

where k is a constant.

Finally, from $x = r \cos\theta$, $y = r \sin\theta$, it follows in the usual way that

$$xy' - yx' = r^2\theta',$$

and hence, by (**) and (*), that

$$2A'(t) = k.$$

Assuming $A(0) = 0$, we conclude that

$$A(t) = 2kt.$$

Thus the radius vector joining the body to the point of origin sweeps out equal areas in equal times (Kepler's law).

Einstein's Use of Infinitesimals

Here is an appropriate place to remark on an intriguing use of infinitesimals in Einstein's celebrated 1905 paper *On the Electrodynamics of Moving Bodies*²⁴, in which the special theory of relativity is first formulated. In deriving the Lorentz transformations from the principle of the constancy of the velocity of light Einstein obtains the following equation for the time coordinate $\tau(x', y, z, t)$ of a moving frame:

$$(i) \quad \frac{1}{2} \left[\tau(0, 0, 0, t) + \tau \left(0, 0, 0, t + \frac{x'}{c-v} + \frac{x'}{c+v} \right) \right] = \tau \left(x', 0, 0, t + \frac{x'}{c-v} \right).$$

He continues:

Hence, if x' be chosen infinitesimally small,

$$(ii) \quad \frac{1}{2} \left(\frac{1}{c-v} + \frac{1}{c+v} \right) \frac{\partial \tau}{\partial t} = \frac{\partial \tau}{\partial x'} + \frac{1}{c-v} \frac{\partial \tau}{\partial t},$$

or

$$\frac{\partial \tau}{\partial x'} + \frac{v}{c^2 - v^2} \frac{\partial \tau}{\partial t} = 0. \text{ }^{25}$$

²⁴ Reprinted in English translation in Einstein et al. (1952). It should be noted, however, that in subsequent presentations of special relativity Einstein avoided the use of infinitesimals

²⁵ Einstein et al. (1952), p. 44.

Now the derivation of equation (ii) from equation (i) can be simply and rigorously carried out in SIA by *choosing* x' to be a *microquantity* ε . For then (i) becomes

$$\frac{1}{2} \left[\tau(0,0,0,t) + \tau \left(0,0,0,t + \varepsilon \left(\frac{1}{c-v} + \frac{1}{c+v} \right) \right) \right] = \tau \left(\varepsilon, 0, 0, t + \frac{\varepsilon}{c-v} \right).$$

From this we get, using equations (1) and (2) above,

$$\tau(0,0,0,t) + \frac{1}{2} \varepsilon \left(\frac{1}{c-v} + \frac{1}{c+v} \right) \frac{\partial \tau}{\partial t} = \tau(0,0,0,t) + \varepsilon \left(\frac{\partial \tau}{\partial x'} + \frac{1}{c-v} \frac{\partial \tau}{\partial t} \right).$$

So

$$\frac{1}{2} \varepsilon \left(\frac{1}{c-v} + \frac{1}{c+v} \right) \frac{\partial \tau}{\partial t} = \varepsilon \left(\frac{\partial \tau}{\partial x'} + \frac{1}{c-v} \frac{\partial \tau}{\partial t} \right),$$

and (ii) follows by microcancellation.

The Lie Bracket

We consider the simple case of vector fields on \mathbf{R} , that is, maps $\mathbf{R} \times \Delta \rightarrow \mathbf{R}$. Any such map X may be expressed in the form

$$X(x, \varepsilon) = x + \varepsilon A(x)$$

for some (unique) map $A: \mathbf{R} \rightarrow \mathbf{R}$. For fixed $x \in \mathbf{R}$ the map $\varepsilon \mapsto \varepsilon A(x): \Delta \rightarrow \mathbf{R}$ is the *field vector* of X at x . The map A is called the *associated map* of X .

The set $V(\mathbf{R})$ of vector fields on \mathbf{R} is a module over the ring \mathbf{R} , with the sum of vector fields and scalar product defined pointwise in the obvious way. How is $V(\mathbf{R})$ turned into an \mathbf{R} -algebra? Given vector fields X, Y with associated maps A, B , we calculate, for microquantities ε, η ,

$$X(Y(x, \eta), \varepsilon) = x + \varepsilon A(x) + \eta B(x) + \varepsilon \eta B(x) \frac{dA}{dx}.$$

so that

$$(*) \quad X(Y(x, \eta), \varepsilon) - Y(X(x, \varepsilon), \eta) = \varepsilon \eta \left[B(x) \frac{dA}{dx} - A(x) \frac{dB}{dx} \right].$$

Now let $[X, Y]$ be the vector field with associated map $B \frac{dA}{dx} - A \frac{dB}{dx}$.

Then from (*) we see that

$$[X, Y](x, \varepsilon \eta) = x + X(Y(x, \eta), \varepsilon) - Y(X(x, \varepsilon), \eta).$$

$[X, Y]$ is the *Lie bracket* of X and Y . With the Lie bracket operation $V(\mathbf{R})$ becomes a Lie algebra over \mathbf{R} .

Given $f: \mathbf{R} \rightarrow \mathbf{R}$ and a vector field X , the *directional derivative* of f along X is the unique map $X(f): \mathbf{R} \rightarrow \mathbf{R}$ satisfying, for all microquantities ε ,

$$f(X(x, \varepsilon)) = f(x) + \varepsilon X(f)(x).$$

If X has associated map A , it is easily checked that $X(f) = A \frac{df}{dx}$. It

follows from this that

$$[X, Y](f) = X(Y(f)) - Y(X(f)).$$

Spacetime Metrics

Spacetime metrics have some arresting properties in SIA. In a spacetime the metric can be written in the form

$$(*) \quad ds^2 = \sum g_{\mu\nu} dx_\mu dx_\nu \quad \mu, \nu = 1, 2, 3, 4.$$

In the classical setting (*) is in fact an abbreviation for an equation involving derivatives and the “differentials” ds and dx_μ are not really quantities at all. What form does this equation take in SIA? Notice that the “differentials” cannot be taken as microquantities since all the squared terms would vanish. But the equation does have a very natural form in terms of microquantities. Here is an informal way of obtaining it.

We think of the dx_μ as being multiples $k_\mu e$ of some small quantity e . Then (*) becomes

$$ds^2 = e^2 \sum g_{\mu\nu} k_\mu k_\nu,$$

so that

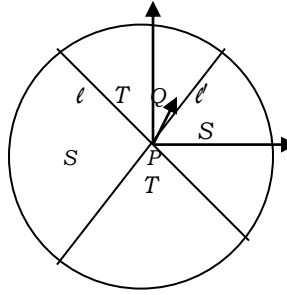
$$ds = e \sqrt{\sum g_{\mu\nu} k_\mu k_\nu}.$$

Now replace e by a microquantity ε . Then we obtain the metric relation in SIA:

$$ds = \varepsilon \sqrt{\sum g_{\mu\nu} k_\mu k_\nu}.$$

This tells us that the “infinitesimal distance” ds between a point P with coordinates (x_1, x_2, x_3, x_4) and an infinitesimally near point Q with coordinates $(x_1 + k_1\varepsilon, x_2 + k_2\varepsilon, x_3 + k_3\varepsilon, x_4 + k_4\varepsilon)$ is $\varepsilon\sqrt{\sum g_{\mu\nu}k_\mu k_\nu}$. Here a curious situation arises. For when the “infinitesimal interval” ds between P and Q is timelike (or lightlike), the quantity $\sum g_{\mu\nu}k_\mu k_\nu$ is nonnegative, so that its square root is a real number. In this case ds may be written as εd , where d is a real number. On the other hand, if ds is spacelike, then $\sum g_{\mu\nu}k_\mu k_\nu$ is negative, so that its square root is imaginary. In this case, then, ds assumes the form $i\varepsilon d$, where d is a real number (and, of course $i = \sqrt{-1}$). On comparing these we see that, if we take ε as the “infinitesimal unit” for measuring infinitesimal timelike distances, then $i\varepsilon$ serves as the “imaginary infinitesimal unit” for measuring infinitesimal spacelike distances.

For purposes of illustration, let us restrict the spacetime to two dimensions (x, t) , and assume that the metric takes the simple form $ds^2 = dt^2 - dx^2$. The infinitesimal light cone at a point P divides the infinitesimal neighbourhood at P into a timelike region T and a spacelike region S bounded by the null lines ℓ and ℓ' respectively (see figure 9). If we take P as origin of coordinates, a typical point Q in this neighbourhood will have coordinates $(a\varepsilon, b\varepsilon)$ with a and b real numbers: if $|b| > |a|$, Q lies in T ; if $a = b$, Q lies on ℓ or ℓ' ; if $|a| < |b|$, Q lies in S . If we write $d = \sqrt{|a^2 - b^2|}$, then in the first case, the infinitesimal distance between P and Q is εd , in the second, it is 0, and in the third it is $i\varepsilon d$.



Minkowski introduced “ ict ” to replace the “ t ” coordinate so as to make the metric of relativistic spacetime positive definite. This was purely a matter of formal convenience, and was later rejected by (general) relativists²⁶. In conventional physics one never works with nilpotent quantities so it is always possible to replace formal imaginaries by their (negative) squares. But spacetime theory in SIA *forces* one to use imaginary units, since, infinitesimally, one can’t “square oneself out of trouble”. This being the case, it would seem that, infinitesimally, the dictum *Farewell to ict*²⁷ needs to be replaced by

Vale “ict”, ave “iε” !

To quote a well-known treatise on the theory of gravitation,

*Another danger in curved spacetime is the temptation to regard ... the tangent space as lying in spacetime itself. This practice can be useful for heuristic purposes, but is incompatible with complete mathematical precision.*²⁸

²⁶ See, for example Box 2.1, *Farewell to “ict”*, of Misner, Thorne and Wheeler (1973).

²⁷ See footnote immediately above.

²⁸ *Op. cit.*, p.205.

The consistency of smooth infinitesimal analysis shows that, on the contrary, yielding to this temptation is compatible with complete mathematical precision: there tangent spaces may indeed be regarded as lying in spacetime itself.

A Speculation

Observe that the microobject Δ is “tiny” in the order-theoretic sense. For, using ε, η as variables ranging over Δ , it is easily seen that that

$$(*) \quad \forall \varepsilon \forall \eta \neg(\varepsilon < \eta \vee \eta < \varepsilon),^{29}$$

whence

$$\forall \varepsilon \forall \eta \varepsilon \leq \eta \wedge \eta \leq \varepsilon.$$

In particular, the members of Δ are all simultaneously ≤ 0 and ≥ 0 , but cannot (because of the nondegeneracy of Δ) be shown to coincide with zero.

In his recent book *Just Six Numbers* the astrophysicist Martin Rees comments on the microstructure of space and time, and the possibility of developing a theory of quantum gravity. In particular he says:

Some theorists are more willing to speculate than others. But even the boldest acknowledge the “Planck scales” as an ultimate barrier. We cannot measure distances smaller than

²⁹ Here is the proof. If microquantities ε, η satisfied $\varepsilon < \eta$, then $0 < \varepsilon - \eta$ so that there is x for which $(\varepsilon - \eta)x = 1$. Squaring both sides gives $1 = (\varepsilon - \eta)^2 x^2 = -2\varepsilon\eta x^2$; squaring both sides of this gives $1 = 0$, a contradiction.

the Planck length [about 10^{19} times smaller than a proton]. We cannot distinguish two events (or even decide which came first) when the time interval between them is less than the Planck time (about 10^{-43} seconds).

On this account, Planck scales seem very similar in certain respects to Δ . In particular, the sentence (*) above seems to be an exact embodiment of the idea that we cannot decide of two “events” in Δ which came first; in fact it makes the stronger assertion that actually neither comes “first”.

Could Δ provide a good model for “Planck scales”? While Δ is unquestionably small enough to play the role, it inhabits a domain in which everything is smooth and continuous, while Planck scales live in the quantum world which, if not outright discrete, is far from being continuous. So if Planck scales could indeed be modelled by microneighbourhoods in SIA, then one might begin to suspect that the quantum microworld, the Planck regime—smaller, in Rees’s words, “than atoms by just as much as atoms are smaller than stars”—is not, like the world of atoms, discrete, but instead continuous like the world of stars. This would be a major victory for the continuous in its long struggle with the discrete.