

Gödel's Dialectica Interpretation

Klaus Frovin Jørgensen

Section for Philosophy and Science Studies, RUC

May 6, 2010

Question to be Answered

Kurt Gödel devised (around 1941) an interpretation of intuitionistic arithmetic into a calculus of functionals.

Together with Gödel's $\neg\neg$ -translation (1933) this was also a consistency proof classical arithmetic.

What are the philosophical conclusions?

The Context of the Interpretation

The 1920s and 1930s:

- The *Grundlagestreit*
- Hilbert's programme (1920s)
- The development of intuitionistic logic (Heyting, Kolmogorov ~ 1930)
- The incompleteness theorems (1931) and the Gödel/Gentzen $\neg\neg$ -translation (1933)
- Church's and Turing's negative answer to the *Entscheidungsproblem* (1935-6).
- Gentzen's proof of consistency of formal number theory (1936)

Heyting's Proof Interpretation

$p : A$

Some of the clauses of Heyting's interpretation are:

- (\rightarrow) $p : A \rightarrow B$ iff p is a construction taking any q such that $q : A$ into $p(q)$ such that $p(q) : B$.
- (\vee) $p : A \vee B$ iff p is a pair (p_0, p_1) , $p_0 \in \{0, 1\}$ and $p_1 : A$ if $p_0 = 0$ and $p_1 : B$ if $p_0 = 1$. $q : A$ into $p(q)$ such that $p(q) : \perp$.
- (\forall) $p : \forall x A(x)$ iff p is a construction taking any t from the intended domain into $p(t)$ such that $p(t) : A(t)$.
- (\exists) $p : \exists x A(x)$ iff p is a pair (p_0, p_1) , where p_0 is an object of the domain and $p_1 : A(p_0)$.

Intuitionistic Propositional Logic

The following is Spector's (1962) formulation:

$$A \rightarrow A, \quad \perp \rightarrow A,$$

$$A \rightarrow A \vee B, \quad B \rightarrow A \vee B,$$

$$A \wedge B \rightarrow A, \quad A \wedge B \rightarrow B,$$

$$\frac{A \quad A \rightarrow B}{B}$$

$$\frac{A \rightarrow B \quad B \rightarrow C}{A \rightarrow C}$$

$$\frac{A \rightarrow B \quad A \rightarrow C}{A \rightarrow B \wedge C}$$

$$\frac{A \wedge B \rightarrow C}{A \rightarrow (B \rightarrow C)}$$

$$\frac{A \rightarrow (B \rightarrow C)}{A \wedge B \rightarrow C}$$

$$\frac{A \rightarrow C \quad B \rightarrow C}{A \vee B \rightarrow C}$$

Tertium Non Datur

Tertium non datur:

$$A \vee \neg A$$

is not sound under the proof interpretation.

Intuitionistic Quantifier Rules

$$\frac{B \rightarrow A(b)}{B \rightarrow \forall x A(x)}$$

$$\forall x A(x) \rightarrow A(t),$$

$$A(t) \rightarrow \exists x A(x),$$

$$\frac{A(b) \rightarrow B}{\exists x A(x) \rightarrow B}$$

Here b is eigenvariable meaning that b is not allowed to occur free in B .

Disjunction Property and Existence Property

- *Existence property*: If $S \vdash \exists x A(x)$ then $S \vdash A(t)$ for a certain term t .
- *Disjunction property*: If $S \vdash A \vee B$, for A, B closed then $S \vdash A$ or $S \vdash B$.

Impredicative Methods (1/2)

A definition is impredicative if it refers to a collection which contains the object to be defined. If one sees definitions as somehow creating or constructing, then circularity is involved. For example the 'least upper bound' of a set is defined to be the 'smallest among the upper bounds'. This is seen, for instance, in the completeness axiom:

Any non-empty subset of real numbers bounded above has a least upper bound.

Example: There exists a real number x such that $x^2 = 2$.

Another example is the Russell set R .

Impredicative Methods (2/2)

Yet an example is the intuitionistic meaning of implication: The proof interpretation). We know $A \rightarrow B$ precisely when we know what may count as a proof of $A \rightarrow B$. A proof of $A \rightarrow B$ transforms any proof of A into a proof of B .

The System Σ (1/3)

The ground type consists of the natural numbers. There are symbols for zero and successor and variables of all types.

If F is an operation of type $\sigma \rightarrow \tau$ this is written as $F^{\sigma \rightarrow \tau}$

Operations in T are defined from combinators (which introduce λ -abstraction)

$$K(x, y) = x \quad S(x, y, z) = x(z)(yz)$$

The combinators give us λ -abstraction $\lambda x.t$ for terms t , with the following equality:

$$(\lambda x.t[x])s = t[s].$$

Moreover, we have primitive recursion

$$\begin{aligned} R(x, y, 0) &= x \\ R(x, y, (z + 1)) &= y(R(x, y, z), z) \end{aligned}$$

The System Σ (2/3)

Quantifier free induction

$$\frac{A(0) \quad A(x^0) \rightarrow A(Sx^0)}{A(x^0)}$$

Substitution:

$$\frac{A(x^\sigma)}{A(t^\sigma)}$$

The System Σ (3/3)

We can do elementary arithmetic in Σ :

$$\begin{aligned}x + y &::= R_0x(\lambda w, u.Sw)y \\ \text{prd} &::= \lambda x.R_00(\lambda w, u.u)x \\ x \dot{\div} y &::= R_0x(\lambda w, u.\text{prd}(w))y \\ |x - y| &::= (x \dot{\div} y) + (y \dot{\div} x) \\ \text{Cond} &::= \lambda x^\sigma, y^\sigma, z^0.R_\sigma x(\lambda v^\sigma, w^0.y)z \\ \text{max} &::= \lambda x^0, y^0.\text{Cond}(y, x, x \dot{\div} y) \\ \text{min} &::= \lambda x^0, y^0.\text{Cond}(x, y, x \dot{\div} y)\end{aligned}$$

The prime formulas of Σ are decidable (proved by induction on the complexity of the terms). Therefore, Σ actually has classical logic which can, moreover, be represented in the system.

Definition of D-translation (1/2)

To each formula A of $\mathcal{L}(\text{HA})$ is now associated its Dialectica translation A^D which is a formula of Σ .

$$A^D \equiv \exists \mathbf{x} \forall \mathbf{y} A_D(\mathbf{x}, \mathbf{y}),$$

where A_D is quantifier free. Intuitively: If A is provable then according to the translation A^D there are \mathbf{x} making A_D 'true' for any \mathbf{y} .

The lengths and types of the fresh variables \mathbf{x} and \mathbf{y} depend only on the logical structure of A .

Definition of D-translation (2/2)

Let A^D have the form $\exists \mathbf{x} \forall \mathbf{y} A_D(\mathbf{x}, \mathbf{y})$ and B^D the form $\exists \mathbf{u} \forall \mathbf{v} B_D(\mathbf{u}, \mathbf{v})$.

Definition.

$$\begin{aligned} A^D &:\equiv A_D \equiv A, \text{ if } A \text{ is prime,} \\ (A \wedge B)^D &:\equiv \exists \mathbf{x}, \mathbf{u} \forall \mathbf{y}, \mathbf{v} (A_D(\mathbf{x}, \mathbf{y}) \wedge B_D(\mathbf{u}, \mathbf{v})), \\ (A \vee B)^D &:\equiv \exists z^0, \mathbf{x}, \mathbf{u} \forall \mathbf{y}, \mathbf{v} ((z = 0 \rightarrow A_D(\mathbf{x}, \mathbf{y})) \wedge \\ &\quad (z \neq 0 \rightarrow B_D(\mathbf{u}, \mathbf{v}))), \\ (\exists z A(z))^D &:\equiv \exists z, \mathbf{x} \forall \mathbf{y} A_D(\mathbf{x}, \mathbf{y}, z), \\ (\forall z A(z))^D &:\equiv \exists \mathbf{X} \forall \mathbf{y}, z A_D(\mathbf{X}z, \mathbf{y}, z), \\ (A \rightarrow B)^D &:\equiv \exists \mathbf{U}, \mathbf{Y} \forall \mathbf{x}, \mathbf{v} (A_D(\mathbf{x}, \mathbf{Y}\mathbf{x}\mathbf{v}) \rightarrow B_D(\mathbf{U}\mathbf{x}, \mathbf{v})). \end{aligned}$$

The Translation of '→'

Assume $\exists x \forall y A_D(x, y) \rightarrow \exists u \forall v B_D(u, v)$. A reasonable reading is:

$$\exists U \forall x (\forall y A_D(x, y) \rightarrow \forall v B_D(Ux, v)). \quad (1)$$

What could a possible interpretation of $\forall x C(x) \rightarrow \forall y D(y)$ be?

Given any counter-example to D we can construct a counter-example to C , i.e. $\exists X \forall y (\neg D(y) \rightarrow \neg C(Xy))$. This implies:

$$\exists U \forall x \exists Y' \forall v (\neg B_D(Ux, v) \rightarrow \neg A_D(x, Y'v)). \quad (2)$$

The quantifier free formulas are stable. Therefore, 'by' the proof interpretation we get:

$$\exists U, Y \forall x, v (A_D(x, Yxv) \rightarrow B_D(Ux, v)).$$

On 'Constructive Meaning'

The following principles are not constructively valid *according* to the proof interpretation by Heyting.

- Different forms of independence-of-premise:

$$(A \rightarrow \exists y B(y)) \rightarrow \exists y (A \rightarrow B(y)),$$

where $y \notin \text{FV}(A)$ and different restrictions on A .

- Markov's principle: $\neg\neg\exists x A_{\text{qf}}(x) \rightarrow \exists x A_{\text{qf}}(x)$.

Gödel's translation changes the intuitionistic meaning:

$$\text{HA}^\omega + \text{IP}_{\forall}^\omega + \text{MP}^\omega + \text{AC} \vdash A \leftrightarrow A^D.$$

This is not necessarily a bad thing.

Theorem [Gödel 1941].

If $\text{HA} \vdash A$, then $\Sigma \vdash A_D(\mathbf{T}, \mathbf{y})$,

where \mathbf{T} is a sequence of terms which can be extracted from a proof of A in HA .

Examples from the Proof of Soundness (1/2)

The proof is by induction on the length of the proof of A in HA.

Case 1. Axiom $A \rightarrow A$. This translates to

$$\exists \mathbf{X}, \mathbf{Y} \forall \mathbf{x}, \mathbf{y} (A_D(\mathbf{x}, \mathbf{Y} \mathbf{x} \mathbf{y}) \rightarrow A_D(\mathbf{X} \mathbf{x}, \mathbf{y})).$$

From this we see that with $\mathbf{T}_1 := \lambda \mathbf{x}. \mathbf{x}$ and $\mathbf{T}_2 := \lambda \mathbf{x}, \mathbf{y}. \mathbf{y}$ we have

$$\Sigma \vdash A_D(\mathbf{x}, \mathbf{T}_2 \mathbf{x} \mathbf{y}) \rightarrow A_D(\mathbf{T}_1 \mathbf{x}, \mathbf{y}).$$

Examples from the Proof of Soundness (2/2)

Case 2. Modus Ponens. Assume as induction hypothesis

- (i) $\Sigma \vdash A_D(\mathbf{T}_1, \mathbf{y})$,
- (ii) $\Sigma \vdash A_D(\mathbf{x}, \mathbf{T}_2 \mathbf{x} \mathbf{v}) \rightarrow B_D(\mathbf{T}_3 \mathbf{x}, \mathbf{v})$,

for given \mathbf{T}_1 , \mathbf{T}_2 , and \mathbf{T}_3 . Find \mathbf{T}_4 such that $\Sigma \vdash B_D(\mathbf{T}_4, \mathbf{v})$. Set \mathbf{x} in (ii) to \mathbf{T}_1 and let \mathbf{y} in (i) be $\mathbf{T}_2 \mathbf{T}_1 \mathbf{v}$. Then use MP (in Σ) to obtain

$$\Sigma \vdash B_D(\mathbf{T}_3 \mathbf{T}_1, \mathbf{v}).$$

Let \mathbf{T}_4 be equal to $\mathbf{T}_3 \mathbf{T}_1$.

Note the similarities to the proof interpretation.

Gödel's Results

In 1941 the following results are mentioned:

- 1 For a certain quantifier free formula $A(x)$ let $C \equiv \neg\forall x(A(x) \vee \neg A(x))$. Then $\text{HA} + C$ is consistent.
- 2 If HA proves $\exists x A(x)$ then Σ proves the *translated* formula $A_D(t)$, for a term t .
- 3 $\neg\neg$ -translation (1933) together with the new interpretation proves consistency of classical arithmetic relative to T .

The interpretation was ultimately published in *Dialectica* in 1958: "Über eine bisher noch nicht benützte Erweiterung des finiten Standpunktes".

Philosophy (1/3)

Weak versus Strong Counterexamples (1/3)

π is transcendental, but we can approximate π arbitrarily well:

$$\pi = \frac{4}{1} - \frac{4}{3} + \frac{4}{5} - \frac{4}{7} + \frac{4}{9} - \frac{4}{11} + \dots$$

Let A be the statement: “There exists one hundred 9s in a row in the decimal expansion of π .”

Weak versus Strong Counterexamples (2/3)

- 1 Start the *computation* of π .
- 2 We start writing a real number a . The first digit is 0 followed by a point:

0.

- 3 Construction of the n -th digit of a :
 - If the decimal expansion of π up to digit number $n - 1$ has not verified A , then the n -th digit of a is 0,
 - Otherwise 1

Is $a = 0$ or $a \neq 0$?

Weak versus Strong Counterexamples (3/3)

Gödel produces with his interpretation a *strong* counterexample. $\neg\forall x(A(x) \vee \neg A(x))$ is demonstrably incompatible with classical logic.

Constructivity is Understood Locally (1/3)

The following principles are not constructively valid *according* to the proof interpretation by Heyting.

- Different forms of independence-of-premise:

$$(A \rightarrow \exists y B(y)) \rightarrow \exists y (A \rightarrow B(y)),$$

where $y \notin \text{FV}(A)$ and different restrictions on A .

- Markov's principle: $\neg\neg\exists x A_{\text{qf}}(x) \rightarrow \exists x A_{\text{qf}}(x)$.

Constructivity is Understood Locally (2/3)

Different interpretations validate different principles:

- Modified realisability validates full extensionality, IP_{ef}^ω and AC.
- Functional interpretation validates MP^ω , IP_{\forall}^ω and AC.

Can there be a better interpretation; one which is more optimal with respect to these principles?

Constructivity is Understood Locally (3/3)

Two incompatible constructive theories:

- The proof interpretation is a global interpretation (a rule of thumb); locally one can accept more if the goal is computable existence.
- The combination of extensionality, Markov's principle and restricted forms of independence-of-premise is a subtle issue.

$$\text{WE-HA}^\omega + \text{IP}_\forall^\omega + \text{MP}^\omega + \text{AC} + \Gamma,$$

Γ is any set of universal true sentences; has existence property, disjunction property and is closed under various rules.

$$\text{E-HA}^\omega + \text{IP}_{\text{ef}}^\omega + \text{AC} + \Gamma,$$

Γ is any set of true \exists -free sentences; has existence property, disjunction property and is closed under different rules, except Markov's rule.

Philosophy (2/3)

On the Existence and Disjunction Properties

Gödel mentions the following result:

- If HA proves $\exists xA(x)$ then Σ proves the *translated* formula $A_D(t)$, for a term t .

But this does not suffice for proving the existence and disjunction properties for HA, as the original formula is not in general intuitionistically provable from the translated. In 1945 Kleene (and Nelson) proved that realisability by numbers can be used for showing that.

Philosophy (3/3):

What should we think of the consistency proof?

Coherence

Gödel's interpretation is paradigmatic with respect to coherence (Σ is now called T):

- PA is interpreted in T
- T is consistent, we can prove strong normalisation, by:
 - Howard's strong computability predicates (uses König's lemma)
 - Tait's method of ascribing ordinals $< \varepsilon_0$ to terms of T
- Fits with Gentzen's partial cut-elimination (which again fits with Schütte's full cut-elimination)
- Tait's proof of termination fits with Gentzen's characterisation of PA as ε_0
- The no-counter-example interpretation can be derived from both the Dialectica interpretation and Gentzen's cut-elimination

Neurath's Ship Welcome Aboard

We cannot eliminate impredicativity. There is no secure Archimedean point—there is no *cogito*.

We give up foundationalism: Instead we can account (I think) for actual mathematics—contemporary as well as historical.