

Logic, Computability and Incompleteness

Gödel's Second Theorem, Löb's
Theorem and the Logic of Provability

Hilbert's Program (again)

- As we saw before, Hilbert's **Formalist Program** for the foundation of mathematics advocated an approach in which all of mathematics is deducible in an **axiomatizable formal theory** where the axioms themselves are **provably consistent**.
- In particular, Hilbert sought an 'internal' and finitary consistency proof for the axioms of elementary number theory.
- **Gödel's First Theorem** can be seen as refuting Hilbert's goal of reducing arithmetic to an 'inventory of provable formulas', while **Gödel's Second Theorem** can be seen as undermining Hilbert's quest for an internal consistency proof for the axioms of arithmetic.

The Unprovability of Consistency

Gödel's First Incompleteness Theorem (roughly): **if** formal arithmetic is consistent, **then** neither S nor $\neg S$ is provable, where S is constructed such that $\vdash S \leftrightarrow \neg \textit{Prov}(\ulcorner S \urcorner)$.

So if arithmetic is consistent then S is unprovable, hence **true**
Is formal arithmetic consistent?

Gödel's Second Incompleteness Theorem (roughly):

if formal arithmetic is consistent,
then it cannot **prove** its own consistency.

A **basic fact** of Classical Logic: if formal arithmetic is a consistent theory, then at least one sentence is unprovable.

The Unprovability of Consistency

So let the **consistency of arithmetic** be expressed
in arithmetic (!) by the sentence:

$$\neg \textit{Prov} (\ulcorner \mathbf{0} = \mathbf{0}' \urcorner) \quad (\underline{\text{con}})$$

Gödel's Second Incompleteness Theorem:

if $\vdash \underline{\text{con}}$ *then* arithmetic is **inconsistent**.

As above, Gödel's First Theorem states that

(#) *if* arithmetic is consistent, *then* it is not provable that S

Since $\underline{\text{con}}$ is a sentence of the object language that expresses the consistency of arithmetic, then a **formalization** of the First Theorem, in terms of (#), would yield:

$$\vdash \underline{\text{con}} \rightarrow \neg \textit{Prov} (\ulcorner S \urcorner) \quad (!)$$

The Unprovability of Consistency

i.e. $\vdash \underline{\text{con}} \rightarrow S$

So *if* this formalization of the First Theorem is indeed provable in formal arithmetic, *then* it follows (on the assumption of consistency) that $\text{not} \vdash \underline{\text{con}}$.

Why? Because $\vdash \underline{\text{con}} \rightarrow S$ entails that

if $\vdash \underline{\text{con}}$ *then* $\vdash S$

And if arithmetic is consistent, then (by 1st Theorem)

$\text{not} \vdash S$, and contraposition on the above yields $\text{not} \vdash \underline{\text{con}}$.

So *if* the conditional $\underline{\text{con}} \rightarrow S$ is a theorem of arithmetic, *then if* arithmetic is consistent, *then* the consistency sentence $\underline{\text{con}}$ cannot be provable in arithmetic.

Gödel's Second Incompleteness Theorem

Thus to **prove** Gödel's Second Incompleteness Theorem **need to show** that $\vdash \text{con} \rightarrow S$, and the underivability of the consequent S will yield the underivability of the antecedent **con**.

proof: will require the **diagonal lemma** and characteristics (i) – (iii) of a **proof predicate**.

As before, the **diagonal lemma** gives

$\vdash S \leftrightarrow \neg \text{Prov}(\ulcorner S \urcorner)$, which yields $\vdash S \rightarrow \neg \text{Prov}(\ulcorner S \urcorner)$

and then $\vdash \neg \neg \text{Prov}(\ulcorner S \urcorner) \rightarrow \neg S$ and finally

(0) $\vdash \text{Prov}(\ulcorner S \urcorner) \rightarrow \neg S$

Recall (i) **if** $\vdash A$, **then** $\vdash \text{Prov}(\ulcorner A \urcorner)$, which applied to (0)

gives $\vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \rightarrow \neg S \urcorner)$

Gödel's Second Incompleteness Theorem

Recall (ii) $\vdash \text{Prov}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}(\ulcorner A \urcorner) \rightarrow \text{Prov}(\ulcorner B \urcorner))$

Which gives $\vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \rightarrow \neg S \urcorner) \rightarrow$
 $(\text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner \neg S \urcorner))$

Now modus ponens gives

(1) $\vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner \neg S \urcorner)$

Similarly, take $\vdash (S \wedge \neg S) \rightarrow \mathbf{o} = \mathbf{o}'$ and apply (i) to get

$\vdash \text{Prov}(\ulcorner (S \wedge \neg S) \rightarrow \mathbf{o} = \mathbf{o}' \urcorner)$ then apply (ii) to get

$\vdash \text{Prov}(\ulcorner (S \wedge \neg S) \rightarrow \mathbf{o} = \mathbf{o}' \urcorner) \rightarrow$

$(\text{Prov}(\ulcorner (S \wedge \neg S) \urcorner) \rightarrow \text{Prov}(\ulcorner \mathbf{o} = \mathbf{o}' \urcorner))$ and MP to get

(\\$) $\vdash \text{Prov}(\ulcorner (S \wedge \neg S) \urcorner) \rightarrow \text{Prov}(\ulcorner \mathbf{o} = \mathbf{o}' \urcorner)$

Gödel's Second Incompleteness Theorem

Note that $\vdash \text{Prov}(\ulcorner (S \wedge \neg S) \urcorner) \leftrightarrow (\text{Prov}(\ulcorner S \urcorner) \wedge \text{Prov}(\ulcorner \neg S \urcorner))$

Substitution of provable equivalents in (\$) yields

$$(2) \quad \vdash (\text{Prov}(\ulcorner S \urcorner) \wedge \text{Prov}(\ulcorner \neg S \urcorner)) \rightarrow \text{Prov}(\ulcorner \mathbf{o} = \mathbf{o}' \urcorner)$$

Recall (iii) $\vdash \text{Prov}(\ulcorner A \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner A \urcorner) \urcorner)$

which gives: (3) $\vdash \text{Prov}(\ulcorner S \urcorner) \rightarrow \underline{\text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \urcorner)}$

Recall (1) $\vdash \underline{\text{Prov}(\ulcorner \text{Prov}(\ulcorner S \urcorner) \urcorner)} \rightarrow \text{Prov}(\ulcorner \neg S \urcorner)$

which in combination with (3) yields

$$(4) \quad \vdash \text{Prov}(\ulcorner S \urcorner) \rightarrow (\text{Prov}(\ulcorner S \urcorner) \wedge \text{Prov}(\ulcorner \neg S \urcorner))$$

With (4) and (2) we get $\vdash \text{Prov}(\ulcorner S \urcorner) \rightarrow \text{Prov}(\ulcorner \mathbf{o} = \mathbf{o}' \urcorner)$

And by contraposition

$$(5) \quad \vdash \neg \text{Prov}(\ulcorner \mathbf{o} = \mathbf{o}' \urcorner) \rightarrow \neg \text{Prov}(\ulcorner S \urcorner).$$

Gödel's Second Incompleteness Theorem

Since con has been defined as the sentence $\neg \text{Prov}(\ulcorner \mathbf{0} = \mathbf{0}' \urcorner)$
and $\vdash S \leftrightarrow \neg \text{Prov}(\ulcorner S \urcorner)$,

Rewriting (5) $\vdash \neg \text{Prov}(\ulcorner \mathbf{0} = \mathbf{0}' \urcorner) \rightarrow \neg \text{Prov}(\ulcorner S \urcorner)$

as $\vdash \text{con} \rightarrow S$ yields the desired result

and if arithmetic is consistent then not $\vdash \text{con}$ \square

This is another incompleteness result, because if arithmetic is consistent then con is *true*,

so if arithmetic is consistent then con is yet another **unprovable truth**, and the wedge between

truth and *provability* that started with the First Theorem is driven even deeper.

Löb's Theorem

We've just seen a 'direct' proof of Gödel's Second Incompleteness Theorem. However, it is also possible to prove this theorem as a corollary of the closely related but more general Löb's Theorem, which is motivated as follows.

Another way to think of provable consistency is in terms of the characteristic

for all sentences A ,

$$(\mathbf{v}) \quad \vdash \textit{Prov} (\ulcorner A \urcorner) \rightarrow A$$

Which 'asserts that' if a sentence is provable, then it is true, so that it's provable in the formal theory that only truths are provable.

Löb's Theorem

Löb's Theorem: if $B(y)$ is a **proof predicate** for some theory T that extends Q , then for any sentence A in the language of T

if $\vdash_T B(\ulcorner A \urcorner) \rightarrow A$ *then* $\vdash_T A$

proof: assume (1) $\vdash_T B(\ulcorner A \urcorner) \rightarrow A$.

Let $D(y)$ be the formula $B(y) \rightarrow A$.

The diagonal lemma guarantees a sentence C such that

$\vdash_T C \leftrightarrow D(\ulcorner C \urcorner)$, i.e.

(2) $\vdash_T C \leftrightarrow (B(\ulcorner C \urcorner) \rightarrow A)$

(1) and (2) in combination with (i) – (iii) yield $\vdash_T A$

(see B&J p. 187 for the detailed steps).

Löb's Theorem

Reformulation of Gödel's Second Incompleteness Theorem:
if $B(y)$ is a **proof predicate** for some theory T that extends Q ,
then **not** $\vdash_T \neg B(\ulcorner o = o' \urcorner)$

proof: suppose $\vdash_T \neg B(\ulcorner o = o' \urcorner)$.

Then $\vdash_T B(\ulcorner o = o' \urcorner) \rightarrow o = o'$ (by prop. logic)

and by Löb's Theorem $\vdash_T o = o'$

But $\vdash_Q \neg o = o'$, and T is inconsistent \square

So **if** T is consistent, **then** it can't prove its own consistency.

The Henkin Sentence

Historically, Löb's Theorem was used to answer a question posed by Henkin with regard to the 'Henkin Sentence' H . Unlike the Gödel sentence S , H 'asserts its own **provability**'. The diagonal lemma guarantees an H such that

$$\vdash H \leftrightarrow \textit{Prov} (\ulcorner H \urcorner)$$

Is H provable (and hence **true**)?

It follows directly from Löb's Theorem that $\vdash H$.

Contrast with the 'truth teller', which is contingent.

Modal Logic of Provability

The defining characteristics of a provability predicate have very clear analogues in modal logic:

(i) $\text{if } \vdash A, \text{ then } \vdash \textit{Prov}(\ulcorner A \urcorner)$

corresponds to the modal inference rule of *Necessitation*

$\text{if } \vdash A, \text{ then } \vdash \Box A$

(ii) $\vdash \textit{Prov}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\textit{Prov}(\ulcorner A \urcorner) \rightarrow \textit{Prov}(\ulcorner B \urcorner))$

corresponds to the **K** axiom schema

$\Box (A \rightarrow C) \rightarrow (\Box A \rightarrow \Box C)$

Modal Logic of Provability

$$(iii) \quad \vdash \textit{Prov} (\ulcorner A \urcorner) \rightarrow \textit{Prov} (\ulcorner \textit{Prov} (\ulcorner A \urcorner) \urcorner)$$

corresponds to the **S4** axiom schema

$$\Box A \rightarrow \Box \Box A$$

$$(v) \quad \vdash \textit{Prov} (\ulcorner A \urcorner) \rightarrow A$$

corresponds to the **T** axiom schema

$$\Box A \rightarrow A$$

The modal theory **S4** is the closure of all the **K**, **S4** and **T** axioms under logical consequence and the rule of *Necessitation*.

Modal Logic of Provability

Hence **S4** is too strong to represent the logic of arithmetical proof, as shown by the Gödel-Löb results.

Instead, need to replace **T** axiom schema with the formalized and then modalized version of Löb's Theorem.

Recall **Löb's Theorem**:

$$\text{if } \vdash \textit{Prov} (\ulcorner A \urcorner) \rightarrow A \text{ then } \vdash A$$

Löb's Theorem formalized in arithmetic:

$$\vdash \textit{Prov} (\ulcorner \textit{Prov} (\ulcorner A \urcorner) \rightarrow A \urcorner) \rightarrow \textit{Prov} (\ulcorner A \urcorner)$$

The **modal** version of the formalization yields the **G axiom schema**:

$$\Box (\Box A \rightarrow A) \rightarrow \Box A$$

Modal Logic of Provability

- The modal theory **G** is the closure of all the **K**, **S4** and **G** axioms under logical consequence and the rule of *Necessitation*.
- Hence the modal theory **G** represents the logic of provability in formal arithmetic.

Montague and Predicate Modal Logic

Suppose the modal concept of **Necessity** is formalized as a 1-place metalinguistic predicate $N(x)$

attaching to **names of formulas**,

rather than as an **operator on formulas**, i.e. \Box .

So the assertion that **it is necessarily the case that Φ**

is formalized as $N(\ulcorner \Phi \urcorner)$ rather than as $\Box \Phi$

Then if the modal theory in question incorporates formal arithmetic and the comparatively weak modal structure of the rule of *Necessitation* (i) and the **T** axiom schema (v), then the theory is **inconsistent**....

Montague and Predicate Modal Logic

proof: the diagonal lemma guarantees a sentence M such that

$$\vdash M \leftrightarrow \neg N(\ulcorner M \urcorner) \quad \text{So ...}$$

- (1) $\vdash M \leftrightarrow \neg N(\ulcorner M \urcorner)$ by diagonal lemma
- (2) $\vdash N(\ulcorner M \urcorner) \rightarrow \neg M$ from (1)
- (3) $\vdash N(\ulcorner M \urcorner) \rightarrow M$ by (v)
- (4) $\vdash \neg N(\ulcorner M \urcorner)$ prop log on (2), (3)
- (5) $\vdash M$ by (1), (4)
- (6) $\vdash N(\ulcorner M \urcorner)$ (i) applied to (5)
- (7) $\vdash N(\ulcorner M \urcorner) \wedge \neg N(\ulcorner M \urcorner)$ from (4), (6)

And the modal theory is inconsistent \square

Leibniz's Law

Leibniz's Law is the principle that the *truth-value of a statement should be preserved under the substitution of co-referential terms*.

It can be seen as a direct corollary of Frege's principle of compositionality:

the semantic value of the whole is a **function of** the semantic values of the relevant parts and their **mode of combination**.

Leibniz's Law holds in all purely extensional contexts.

However, the law can **fail** in propositional attitude contexts such as *knowledge* and *belief*.

Failure of Leibniz's Law

For example, consider the following:

- (i) Aristotle knew **that** 9 > 7 (true)
- (ii) 9 = the number of planets (true)
- (iii) Aristotle knew **that** the number of planets > 7 (false)

Or consider an example using belief

- (i) Frank believes **that** gold is valuable (true)
- (ii) Gold is the element with atomic number 79 (true)
- (iii) Frank believes **that** the element with atomic number 79 is valuable (*de dicto* false, if Frank is unaware of (ii))

Failure of Leibniz's Law in Metamathematics

Let S be the Gödel sentence, so that (again)

$$\vdash S \leftrightarrow \neg \textit{Prov} (\ulcorner S \urcorner).$$

If arithmetic is consistent then

S is true and unprovable, while

$\neg S$ is false and unprovable.

Suppose $\ulcorner S \urcorner = \mathbf{a}$.

Then S is equivalent to the sentence $\neg \textit{Prov} (\mathbf{a})$.

Let d be the definite description (using Russell's variable-binding, term-forming iota operator (ιx),

to be read as 'the x such that...')

$$d = (\iota x) ((\neg \textit{Prov} (\mathbf{a}) \rightarrow x = \mathbf{a}) \wedge (\textit{Prov} (\mathbf{a}) \rightarrow x = \ulcorner \mathbf{0} = \mathbf{0}' \urcorner))$$

Failure of Leibniz's Law in Metamathematics

If arithmetic is consistent then $d = a$.

But this truth is not provable in arithmetic,
because it would require proving $\neg \textit{Prov}(a)$.

Hence, we have the situation:

- (1) It is provable in arithmetic that $a = a$ (true)
- (2) $a = d$ (true)
- (3) It is provable in arithmetic that $a = d$ (false)

Where (3) is derived from (1) by the substitution of co-referential terms given in (2).

So the context 'It is provable in arithmetic that ...'
violates Leibniz's Law, if arithmetic is consistent.