

16 Proving the Second Incompleteness Theorem

In his original 1931 paper, Gödel states a version of the Second Theorem. But he doesn't spell out a full proof. He simply claims that the reasoning for the First Theorem is so elementary that a strong enough theory must be able to replicate the reasoning and prove the Formalized First Theorem; and then he notes that this implies the Second Theorem.¹

The hard work of taking a strong enough T and checking that it really does prove the Formalized First Theorem was first done for a particular case by David Hilbert and Paul Bernays in their *Grundlagen der Mathematik* of 1939. The details of their proof are – the story goes – due to Bernays, who had discussed it with Gödel during a transatlantic voyage.

Now, Hilbert and Bernays helpfully isolated what we can call *derivability conditions* on the predicate Prov_T , conditions whose satisfaction is indeed enough for a theory T to prove the Formalized First Theorem. Later, Martin H. Löb gave a rather neater version of these conditions, giving us the so-called *HBL conditions* have become standard and which we invoke in this chapter.

16.1 Sharpening the Second Theorem

Recall: $\text{Prov}_T(y)$ abbreviates $\exists v \text{Prf}_T(v, y)$; and Prf_T is a Σ_1 expression. So arguing about provability inside T will involve establishing some general claims involving Σ_1 expressions. And how do we prove general arithmetical claims? Using induction is the default method.

It therefore looks a good bet that we will get $T \vdash \text{Con}_T \rightarrow \neg \text{Prov}_T(\overline{\Gamma \text{G}_T \neg})$ if T at least has Σ_1 -induction – meaning that T 's axioms include (the universal closures of) all instances of the first-order Induction Schema where the induction predicate φ is no more complex than Σ_1 .

In §7.5 we introduced $|\Sigma_1$ as the standard label for the result of adding that amount of induction to \mathbb{Q} . So it is a plausible conjecture that we can in fact show the following, a fleshed-out version of our skeletal Theorem 57:

¹I am told that Gödel's shorthand notebooks at the time suggest that he in fact hadn't then worked out a detailed proof of the first step.

16 Proving the Second Incompleteness Theorem

Theorem 60. *If T is p.r. axiomatized and contains $\mathbf{I}\Sigma_1$, then T proves the Formalized First Theorem, i.e. $T \vdash \text{Con}_T \rightarrow \neg \text{Prov}_T(\overline{\text{G}_T})$*

And if that is right, by the argument of §15.3 we get an improved version of Theorem 58.

Theorem 61. *If T is consistent, p.r. axiomatized and contains $\mathbf{I}\Sigma_1$, then $T \not\vdash \text{Con}_T$.*

So how can we prove Theorem 60?

16.2 The box notation

To improve readability, we introduce some standard notation. We will henceforth abbreviate $\text{Prov}_T(\overline{\varphi})$ simply by $\Box_T \varphi$. This is perhaps a bit naughty, as the new notation hides away some Gödel-numbering. But the notation is standard and is safe enough if we keep our wits about us.²

So in particular, $\neg \text{Prov}_T(\overline{\text{G}_T})$ can be abbreviated $\neg \Box_T \text{G}_T$. Thus in our new notation, the Formalized First Theorem is $\text{Con}_T \rightarrow \neg \Box_T \text{G}_T$. Moreover, Con_T can now alternatively be abbreviated as $\neg \Box_T \perp$.

To reduce clutter, we will also very often drop the explicit subscript T from the box symbol and elsewhere, and let context supply it.

16.3 Proving the Formalized First Theorem

First a standard definition: we will say (dropping subscripts)

Defn. 58. *The derivability conditions hold in T if and only if, for any T -sentences φ, ψ ,*

1. *If $T \vdash \varphi$, then $T \vdash \Box \varphi$;*
2. *$T \vdash \Box(\varphi \rightarrow \psi) \rightarrow (\Box \varphi \rightarrow \Box \psi)$;*
3. *$T \vdash \Box \varphi \rightarrow \Box \Box \varphi$.*

Here, (1) tells us that if T can prove φ , then (via the relevant Gödel-numbering which enables T to code up claims about provability-in- T), T ‘knows’ it can prove φ . (2) tells us that T is aware of modus ponens; if T knows it can prove $(\varphi \rightarrow \psi)$, and knows it can prove φ , then T knows it can prove ψ too. And (3) tells us that if T knows it can prove φ , it can code up that proof so it can prove that φ is provable. We might reasonably expect that these conditions will hold for a strong enough theory T .

And now we can indeed show that

²If you are familiar with modal logic, then you will immediately recognize the conventional symbol for the necessity operator. And the parallels and differences between “‘ $1 + 1 = 2$ ’ is provable (in T)” and ‘It is necessarily true that $1 + 1 = 2$ ’ are highly suggestive. These parallels and differences are the topic of ‘provability logic’, the subject of a contemporary classic, Boolos’s *The Logic of Provability*.

Theorem 62. *If T is p.r. axiomatized and contains $\text{I}\Sigma_1$, then the derivability conditions hold for T .*

However, demonstrating this is a seriously tedious task, and I certainly *don't* propose to do hack through the annoying technical details. In these notes we will just take this technical theorem as given.

We also can show, considerably more easily, that

Theorem 63. *If T is p.r. axiomatized, contains Q , and the derivability conditions hold for T , then T proves the Formalized First Theorem.*

Proof. This is just a mildly fun exercise in box-pushing, with the target of showing $T \vdash \text{Con}_T \rightarrow \neg \Box_T \text{G}_T$.

First, since T is p.r. axiomatized and contains Q , Theorem 50 holds. So, in our new symbolism, we have $T \vdash \text{G} \leftrightarrow \neg \Box \text{G}$.

Second, since T contains Q and standard logic, we have

$$T \vdash \neg \varphi \rightarrow (\varphi \rightarrow \perp).$$

(whichever way we read the absurdity constant – see Defn. 56). Given this and the derivability condition (C1), we get

$$T \vdash \Box(\neg \varphi \rightarrow (\varphi \rightarrow \perp)).$$

So given the derivability condition (C2) and using modus ponens, it follows that for any φ

$$\text{A. } T \vdash \Box \neg \varphi \rightarrow \Box(\varphi \rightarrow \perp).$$

We now argue as follows, dropping the subscript T s:

- | | |
|---|--|
| 1. $T \vdash \text{G} \rightarrow \neg \Box \text{G}$ | Half of Theorem 50 |
| 2. $T \vdash \Box(\text{G} \rightarrow \neg \Box \text{G})$ | From 1, given C1 |
| 3. $T \vdash \Box \text{G} \rightarrow \Box \neg \Box \text{G}$ | From 2, using C2 |
| 4. $T \vdash \Box \neg \Box \text{G} \rightarrow \Box(\Box \text{G} \rightarrow \perp)$ | Instance of A |
| 5. $T \vdash \Box \text{G} \rightarrow \Box(\Box \text{G} \rightarrow \perp)$ | From 3 and 4 |
| 6. $T \vdash \Box \text{G} \rightarrow (\Box \Box \text{G} \rightarrow \Box \perp)$ | From 5, using C2 and logic |
| 7. $T \vdash \Box \text{G} \rightarrow \Box \Box \text{G}$ | Instance of C3 |
| 8. $T \vdash \Box \text{G} \rightarrow \Box \perp$ | From 6 and 7 |
| 9. $T \vdash \neg \Box \perp \rightarrow \neg \Box \text{G}$ | Contraposing |
| 10. $T \vdash \text{Con} \rightarrow \neg \Box \text{G}$ | Definition of Con ☒ |

Now we put the last two theorems together. If T is p.r. axiomatized and contains $\text{I}\Sigma_1$, the derivability conditions will hold (by Theorem 62) and it contains Q . Hence (by Theorem 63) if T is p.r. axiomatized and contains $\text{I}\Sigma_1$, it proves the Formalized First Theorem.

Which gives us our target Theorem 60.

16.4 The equivalence of Con_T with G_T

We'll assume throughout this section that T is a p.r. axiomatized theory such that the derivability conditions hold. And first, we will show that in that case, G_T and Con_T are provably equivalent in T .

We have just shown in the previous Theorem that $T \vdash \text{Con} \rightarrow \neg \Box G$ (suppressing subscripts!). We can now prove the converse, $T \vdash \neg \Box G \rightarrow \text{Con}$. In fact, we have a more general result:

Theorem 64. *For any sentence φ , $T \vdash \neg \Box \varphi \rightarrow \text{Con}$.*

Proof. We argue as follows:

- | | |
|---|--|
| 1. $T \vdash \perp \rightarrow \varphi$ | Logic! |
| 2. $T \vdash \Box(\perp \rightarrow \varphi)$ | From 1, given C1 |
| 3. $T \vdash \Box \perp \rightarrow \Box \varphi$ | From 2, given C2 |
| 4. $T \vdash \neg \Box \varphi \rightarrow \neg \Box \perp$ | Contraposing |
| 5. $T \vdash \neg \Box \varphi \rightarrow \text{Con}$ | Definition of Con ⊠ |

This little theorem has a rather remarkable corollary. Since T can't prove Con , the theorem tells us that T doesn't entail $\neg \Box \varphi$ for any φ at all. Hence T doesn't ever 'know' that it can't prove φ , even when it can't! In sum, suppose that T satisfies the now familiar conditions: by (C1), T knows all about what it *can* prove; but we have just shown that it knows nothing about what it *can't* prove.

Now, as a particular instance of our last theorem, we have $T \vdash \neg \Box G \rightarrow \text{Con}$. So putting that together with Theorem 63, we have $T \vdash \text{Con} \leftrightarrow \neg \Box G$. And now combine *that* with Theorem 50 which tells us that $T \vdash G \leftrightarrow \neg \Box G$, and lo and behold we've shown

Theorem 65. *If T is p.r. axiomatized and contains $\text{I}\Sigma_1$, then $T \vdash \text{Con} \leftrightarrow G$.*

This means that, not only do we have $T \not\vdash \text{Con}$, we also have (assuming in addition that T is ω -consistent) $T \not\vdash \neg \text{Con}$. In other words, Con is formally undecidable by T .

But Con is *not* self-referential. That observation should scotch any lingering suspicion that undecidable sentences provided by Gödelian arguments are all tainted by potentially paradoxical self reference.

We next prove

Theorem 66. $T \vdash \text{Con} \leftrightarrow \neg \Box \text{Con}$.

Proof. The direction from right to left is an instance of Theorem 64. For the other direction, note that

- | | |
|--|------------------|
| 1. $T \vdash \text{Con} \rightarrow G$ | Already proved |
| 2. $T \vdash \Box(\text{Con} \rightarrow G)$ | From 1, given C1 |
| 3. $T \vdash \Box \text{Con} \rightarrow \Box G$ | From 2, given C2 |

- | | | |
|--|---------------|---|
| 4. $T \vdash \neg \Box G \rightarrow \neg \Box \text{Con}$ | Contraposing | |
| 5. $T \vdash \text{Con} \rightarrow \neg \Box \text{Con}$ | Using Thm. 63 | ☒ |

So: this shows that Con (like G) is also a fixed point of the negated provability predicate (see again Defn. 52 and §13.3).

As we have noted before, some authors refer to any fixed point of the negated provability predicate as a Gödel sentence. Fine. That’s one way of using the jargon. But if you adopt the broad usage, you must be careful with your informal commentary. For example, not all Gödel sentences in the broad sense indirectly ‘say’ *I am unprovable*: Con is a case in point.

16.5 Consistent theories that ‘prove’ their own inconsistency

An amusing and instructive afterthought to finish this chapter.

An ω -consistent (and so consistent) T can’t prove $\neg \text{Con}_T$, as we’ve just noted. By contrast, a consistent but ω -inconsistent T can have $\neg \text{Con}_T$ as a theorem. (Being ω -inconsistent, T won’t be sound; that’s how it can have a false theorem!)

The proof is actually straightforward, once we note a simple lemma. Suppose S and R are two p.r. axiomatized theories, which share a deductive logic; and suppose every axiom of the simpler theory S is also an axiom of the richer theory R . It is then a trivial logical truth that, if the richer theory R is consistent, then the simpler theory S must be consistent too. And the arithmetical claim that encodes this fact can also be formally proved:

Theorem 67. *Under the given conditions, with theory R extending theory S , $R \vdash \text{Con}_R \rightarrow \text{Con}_S$.*

So now take our simpler theory S to be PA . Take the richer theory R to be PA augmented by the extra axiom $\neg G_{\text{PA}}$. We have met this richer theory briefly before: by theorem 43, R is consistent but ω -inconsistent. Since R trivially proves $\neg G_{\text{PA}}$, and PA (and hence R) proves G_{PA} is equivalent to Con_{PA} , R proves $\neg \text{Con}_{\text{PA}}$. So – using our last theorem and modus tollens! – R proves $\neg \text{Con}_R$.

Summing that up,

Theorem 68. *Assuming PA is consistent, the theory $\text{PA} + \neg G_{\text{PA}}$ is a consistent theory which ‘proves’ itself inconsistent (i.e. it entails the negation of its own consistency sentence).*

And since R proves $\neg \text{Con}_R$,

Theorem 69. *There can be a consistent theory R is such that the theory $R + \text{Con}_R$ is inconsistent.*

What are we to make of these apparent absurdities? Well, giving the language of R , i.e. $\text{PA} + \neg G_{\text{PA}}$, its standard arithmetical interpretation, the theory is unsound (since the Gödel sentence G_{PA} is true). So we shouldn’t trust what R says about anything, especially about its own inconsistency when we derive

16 Proving the Second Incompleteness Theorem

$\neg\text{Con}_R$ from it. R doesn't really *prove* (in the ordinary sense of establish-as-true) its own inconsistency, since we don't accept the theory as correct on the standard interpretation! That's why we used scare quotes in stating Theorem 68.