

21 Using the Diagonalization Lemma

(c) But now we want to show that we don't need the assumption of soundness: consistency is enough. To show this, we first prove the following general result, which is the analogue of Theorem 21.1:

Theorem 21.2 *Let T be a nice theory, and let γ be any fixed point for $\neg\text{RProv}_T(x)$. Then $T \not\vdash \gamma$ and $T \not\vdash \neg\gamma$.*

Proof for first half Suppose γ is any theorem. Then – dropping subscripts for readability – for some m , $\text{Prf}(m, \ulcorner \gamma \urcorner)$. Since Prf captures Prf , $T \vdash \text{Prf}(\bar{m}, \ulcorner \gamma \urcorner)$.

Also, since T is consistent, $\neg\gamma$ is unprovable, so for all n , not- $\overline{\text{Prf}}(\bar{n}, \ulcorner \gamma \urcorner)$. Since $\overline{\text{Prf}}$ captures $\overline{\text{Prf}}$, then for each $n \leq m$ in particular, $T \vdash \neg\overline{\text{Prf}}(\bar{n}, \ulcorner \gamma \urcorner)$. Using the result (O4) of Section 9.4, that shows $T \vdash (\forall w \leq \bar{m}) \neg\overline{\text{Prf}}(w, \ulcorner \gamma \urcorner)$.

Putting these results together, $T \vdash \text{Prf}(\bar{m}, \ulcorner \gamma \urcorner) \wedge (\forall w \leq \bar{m}) \neg\overline{\text{Prf}}(w, \ulcorner \gamma \urcorner)$. So, existentially quantifying, $T \vdash \text{RProv}(\ulcorner \gamma \urcorner)$.

But now suppose that γ is indeed a fixed point for $\neg\text{RProv}(x)$, i.e. $T \vdash \gamma \leftrightarrow \neg\text{RProv}(\ulcorner \gamma \urcorner)$. Then if γ is provable, we'd also have $T \vdash \neg\text{RProv}(\ulcorner \gamma \urcorner)$. Contradiction. So a fixed point γ is not provable: $T \not\vdash \gamma$. \square

Proof for second half Now suppose $\neg\gamma$ is a theorem, for some γ . Then for some m , $\overline{\text{Prf}}(\bar{m}, \ulcorner \gamma \urcorner)$, so $T \vdash \overline{\text{Prf}}(\bar{m}, \ulcorner \gamma \urcorner)$.

Also, since T is consistent, γ is unprovable, so for all n , not- $\text{Prf}(n, \ulcorner \gamma \urcorner)$. Hence, by a parallel argument to before, $T \vdash (\forall v \leq \bar{m}) \neg\text{Prf}(v, \ulcorner \gamma \urcorner)$. Elementary manipulation gives $T \vdash \forall v (\text{Prf}(v, \ulcorner \gamma \urcorner) \rightarrow \neg v \leq \bar{m})$. Now appeal to (O8) of Section 9.4, and that gives $T \vdash \forall v (\text{Prf}(v, \ulcorner \gamma \urcorner) \rightarrow \bar{m} \leq v)$.

Combining these two results, it immediately follows that $T \vdash \forall v (\text{Prf}(v, \ulcorner \gamma \urcorner) \rightarrow (\bar{m} \leq v \wedge \overline{\text{Prf}}(\bar{m}, \ulcorner \gamma \urcorner)))$. That implies $T \vdash \forall v (\text{Prf}(v, \ulcorner \gamma \urcorner) \rightarrow (\exists w \leq v) \overline{\text{Prf}}(w, \ulcorner \gamma \urcorner))$. So given our definition, $T \vdash \neg\text{RProv}(\ulcorner \gamma \urcorner)$.

Suppose again that γ is a fixed point for $\neg\text{RProv}(x)$, i.e. $T \vdash \gamma \leftrightarrow \neg\text{RProv}(\ulcorner \gamma \urcorner)$. Then if $\neg\gamma$ is provable, we'd also have $T \vdash \text{RProv}(\ulcorner \gamma \urcorner)$. Contradiction. So if γ is a fixed point, $\neg\gamma$ is not provable: $T \not\vdash \neg\gamma$. \square

(d) So we now know that any fixed point for $\neg\text{RProv}_T$ must be formally undecidable in T . But the Diagonalization Lemma has already told us that there has to be such a fixed point R_T . Hence, assuming no more than T 's niceness, it follows that T is negation-incomplete.

Which is almost what we wanted to show. But not quite. For recall our official statement of the Gödel-Rosser Theorem:

Theorem 19.6 *If T is a nice theory, then there is an L_A -sentence φ of Goldbach type such that neither $T \vdash \varphi$ nor $T \vdash \neg\varphi$.*

This says not just that a nice theory T has an undecidable sentence, but that it has a Π_1 undecidable sentence. And how do we show *that*?

This time it isn't enough simply to appeal to the corollary of Theorem 20.4, i.e. to the principle that Π_1 predicates have Π_1 fixed points. For $\neg\text{RProv}(x)_T$

isn't Π_1 (or at least, not evidently so),⁵ so we can't conclude that its fixed point R_T is Π_1 . Hence we are going to have to do a bit more work to demonstrate the full-strength Gödel-Rosser Theorem.

Proof Let's look at the proof of the previous theorem again, and generalize the leading idea.

Suppose, then, that instead of using the two-place predicates Prf and $\overline{\text{Prf}}$ we use any other pair of two-place predicates P and \overline{P} which respectively “enumerate” the positive and negative T -theorems, i.e. satisfy the following conditions:

1. if $T \vdash \gamma$, then for some m , $T \vdash P(\overline{m}, \ulcorner \gamma \urcorner)$.
2. if $T \not\vdash \gamma$, then for all n , $T \vdash \neg P(\overline{n}, \ulcorner \gamma \urcorner)$.
3. if $T \vdash \neg \gamma$, then for some m , $T \vdash \overline{P}(\overline{m}, \ulcorner \gamma \urcorner)$.
4. if $T \not\vdash \neg \gamma$, then for all n , $T \vdash \neg \overline{P}(\overline{n}, \ulcorner \gamma \urcorner)$.

Now define $\text{RP}_T(x) =_{\text{def}} \exists v(P(v, x) \wedge (\forall w \leq v) \neg \overline{P}(w, x))$. This gives us another Rosser-style predicate, and the argument will go through *exactly* as before: for a nice theory T , any fixed point of $\neg \text{RP}_T(x)$ will be undecidable.

This tells what we need to look for. Suppose we can find predicates P and \overline{P} which satisfy our four “enumeration” conditions, but which are Δ_0 (i.e. lack unbounded quantifiers). Then the corresponding $\text{RP}_T(x)$ will evidently be Σ_1 : so its negation $\neg \text{RP}_T(x)$ *will* be Π_1 and will indeed have Π_1 undecidable fixed points.

It just remains, then, to find a suitable pair of Δ_0 predicates P and \overline{P} . Well, consider the Σ_1 formula $\text{Prov}_T(x) =_{\text{def}} \exists v \text{Prf}(v, x)$. That expresses the property Prov_T , i.e. the property of Gödel-numbering a T -theorem (see Section 20.1). Since it is Σ_1 , $\text{Prov}_T(x)$ is logically equivalent to a wff with a bunch of initial existential quantifiers followed by a Δ_0 wff. And we can now apply the same trick we invoked in proving Theorem 10.1 to get a wff that expresses the same property Prov_T but which starts with just a *single* existential quantifier, i.e. has the form $\exists u P(u, x)$ where P is Δ_0 .

But note that when γ is a theorem, $\exists u P(u, \ulcorner \gamma \urcorner)$ is true, so for some m , $P(\overline{m}, \ulcorner \gamma \urcorner)$ is true. So, being nice and hence Δ_0 -complete, T proves that last wff. And if γ isn't a theorem, $\exists u P(u, \ulcorner \gamma \urcorner)$ is false, so for every n , $P(\overline{n}, \ulcorner \gamma \urcorner)$ is false, so each $\neg P(\overline{n}, \ulcorner \gamma \urcorner)$ is true. Being Δ_0 -complete, T proves all those latter wffs too.

Hence P is Δ_0 and satisfies the “enumerating” conditions (1) and (2). We can similarly construct a Δ_0 wff \overline{P} from $\exists v \overline{\text{Prf}}(v, x)$. So we are done. \square

Phew!

⁵Why? Well, $\text{RProv}(x)_T$ is defined as $\exists v(\text{Prf}_T(v, x) \wedge (\forall w \leq v) \neg \overline{\text{Prf}}_T(w, x))$, and its component wff $\neg \overline{\text{Prf}}_T$ is Π_1 . So, after the initial existential quantifier, $\text{RProv}(x)_T$ in effect has an unbounded *universal* quantifier buried inside. Hence $\text{RProv}(x)_T$ isn't strictly Σ_1 : and it isn't evidently logically equivalent to a strictly Σ_1 wff either. So its negation isn't evidently Π_1 .