

## inp.1 The Undefinability of Truth

inc:inp:tar:  
sec

The notion of *definability* depends on having a formal semantics for the language of arithmetic. We have described a set of formulas and sentences in the language of arithmetic. The “intended interpretation” is to read such sentences as making assertions about the natural numbers, and such an assertion can be true or false. Let  $\mathfrak{N}$  be the **structure** with domain  $\mathbb{N}$  and the standard interpretation for the symbols in the language of arithmetic. Then  $\mathfrak{N} \models \varphi$  means “ $\varphi$  is true in the standard interpretation.”

**Definition inp.1.** A relation  $R(x_1, \dots, x_k)$  of natural numbers is *definable* in  $\mathfrak{N}$  if and only if there is a formula  $\varphi(x_1, \dots, x_k)$  in the language of arithmetic such that for every  $n_1, \dots, n_k$ ,  $R(n_1, \dots, n_k)$  if and only if  $\mathfrak{N} \models \varphi(\bar{n}_1, \dots, \bar{n}_k)$ .

Put differently, a relation is definable in  $\mathfrak{N}$  if and only if it is representable in the theory **TA**, where **TA** =  $\{\varphi : \mathfrak{N} \models \varphi\}$  is the set of true sentences of arithmetic. (If this is not immediately clear to you, you should go back and check the definitions and convince yourself that this is the case.)

**Lemma inp.2.** *Every computable relation is definable in  $\mathfrak{N}$ .*

*Proof.* It is easy to check that the formula representing a relation in **Q** defines the same relation in  $\mathfrak{N}$ .  $\square$

Now one can ask, is the converse also true? That is, is every relation definable in  $\mathfrak{N}$  computable? The answer is no. For example:

**Lemma inp.3.** *The halting relation is definable in  $\mathfrak{N}$ .*

*Proof.* Let  $H$  be the halting relation, i.e.,

$$H = \{\langle e, x \rangle : \exists s T(e, x, s)\}.$$

Let  $\theta_T$  define  $T$  in  $\mathfrak{N}$ . Then

$$H = \{\langle e, x \rangle : \mathfrak{N} \models \exists s \theta_T(\bar{e}, \bar{x}, s)\},$$

so  $\exists s \theta_T(z, x, s)$  defines  $H$  in  $\mathfrak{N}$ .  $\square$

**Problem inp.1.** Show that  $Q(n) \Leftrightarrow n \in \{\# \varphi \# : \mathbf{Q} \vdash \varphi\}$  is definable in arithmetic.

What about **TA** itself? Is it definable in arithmetic? That is: is the set  $\{\# \varphi \# : \mathfrak{N} \models \varphi\}$  definable in arithmetic? Tarski’s theorem answers this in the negative.

inc:inp:tar:  
thm:tarski

**Theorem inp.4.** *The set of true statements of arithmetic is not definable in arithmetic.*

*Proof.* Suppose  $\theta(x)$  defined it. By the fixed-point lemma, there is a formula  $\varphi$  such that  $\mathbf{Q}$  proves  $\varphi \leftrightarrow \neg\theta(\ulcorner\varphi\urcorner)$ , and hence  $\mathfrak{N} \models \varphi \leftrightarrow \neg\theta(\ulcorner\varphi\urcorner)$ . But then  $\mathfrak{N} \models \varphi$  if and only if  $\mathfrak{N} \models \neg\theta(\ulcorner\varphi\urcorner)$ , which contradicts the fact that  $\theta(y)$  is supposed to define the set of true statements of arithmetic.  $\square$

Tarski applied this analysis to a more general philosophical notion of truth. Given any language  $L$ , Tarski argued that an adequate notion of truth for  $L$  would have to satisfy, for each sentence  $X$ ,

‘ $X$ ’ is true if and only if  $X$ .

Tarski’s oft-quoted example, for English, is the sentence

‘Snow is white’ is true if and only if snow is white.

However, for any language strong enough to represent the diagonal function, and any linguistic predicate  $T(x)$ , we can construct a sentence  $X$  satisfying “ $X$  if and only if not  $T(\ulcorner X \urcorner)$ .” Given that we do not want a truth predicate to declare some sentences to be both true and false, Tarski concluded that one cannot specify a truth predicate for all sentences in a language without, somehow, stepping outside the bounds of the language. In other words, a truth predicate for a language cannot be defined in the language itself.

## Photo Credits

## Bibliography