

# Notes on Undecidability

## Foundations of Computer Science

Spring 2017

Remember that a language  $A \subseteq \Sigma^*$  is decidable (or Turing-decidable when we want to be more precise) when there is a *total* Turing machine  $M$  that accepts  $A$ . By the Church-Turing thesis, we can say that a decision problem  $d : \Sigma^* \rightarrow \{1, 0\}$  is computable when  $\{u \in \Sigma^* \mid f(u) = 1\}$  is decidable.

A language is said to be *undecidable* when it is not decidable. A simple counting argument shows that there has to be at least one undecidable language. Intuitively, there are more languages (over alphabet  $\Sigma$ ) than there are Turing machines (over alphabet  $\Sigma$ ).

Because there are infinitely many languages and infinitely many Turing machines, making such a statement precise requires a definition of “size” for infinite sets.

**Definition:** Two sets  $A$  and  $B$  are *equipollent* (have the same size), written  $A \approx B$ , if you can match every element of  $A$  with a distinct element of  $B$ , and vice versa. Formally,  $A \approx B$  when there exists a *one-to-one and onto* function (a.k.a. a *bijection*)  $f : A \rightarrow B$ .

**Definition:** Set  $A$  is “no bigger than” set  $B$ , written  $A \preceq B$ , if you can match every element of  $A$  with a distinct element of  $B$ . Formally,  $A \preceq B$  if there  $A \approx C$  for some  $C \subseteq B$ , or equivalently if there exists a one-to-one function  $f : A \rightarrow B$ .

**Definition:**  $A \prec B$  if  $A \preceq B$  but  $A \not\approx B$ .

**Properties:**

- (i) Both  $\approx$  and  $\preceq$  are transitive.
- (ii) (Cantor-Bernstein) If  $A \preceq B$  and  $B \preceq A$  then  $A \approx B$ .
- (iii) If  $A \prec B$ , then there is no onto function  $f : A \rightarrow B$ .

Property (i) is pretty easy to show. Property (ii) is much more challenging. Property (iii) can be proved from property (ii).

If  $\mathbb{N}$  is the set of natural numbers and  $2\mathbb{N}$  the set of even natural numbers, then the bijection  $f : \mathbb{N} \rightarrow 2\mathbb{N}$  given by  $f(n) = 2n$  shows that  $\mathbb{N} \approx 2\mathbb{N}$ . Similarly, you can show that  $\mathbb{N} \approx 2\mathbb{N} + 1$ , where  $2\mathbb{N} + 1$  is the set of odd natural numbers.

It is easy to show that  $\mathbb{N} \preceq \mathbb{Q}$ . To show that  $\mathbb{Q} \preceq \mathbb{N}$ , imagine putting all rational numbers in an (infinite) array with rational number  $i/j$  at the cell in column  $i$  and row  $j$ . We can now associate a natural number with each rational number traversing the array in diagonal bands, where band  $k$  lists all rational number of the form  $i/j$  with  $i + j = k$ . We treat zero specially. Formally, we can define  $f : \mathbb{Q} \rightarrow \mathbb{N}$  by taking

$$\begin{aligned} f(0) &= 0 \\ f(n/m) &= \Delta(n + m - 2) + m - 1 \end{aligned}$$

where  $\Delta(k)$  is the  $k$ th triangular number, defined by  $\Delta(0) = 0$  and  $\Delta(n + 1) = \Delta(n) + n$ . the first triangular numbers are 0, 1, 3, 6, 10, 15, 21, ... . Function  $f$  is one-to-one, so  $\mathbb{Q} \preceq \mathbb{N}$ . By property (ii) above,  $\mathbb{N} \approx \mathbb{Q}$ .

Not every infinite set is equipollent to  $\mathbb{N}$ . A classic argument can be used to show that  $\mathbb{N} \prec \mathbb{R}$ , the set of real numbers.

Rather than show the latter, though, here's a result that will be more useful for us.

**Cantor's Theorem:** For any set  $A$ , we have  $A \prec 2^A$ , where  $2^A$  is the set of subsets of  $A$ .

Notation  $\wp(A)$  is also used for  $2^A$ .

*Proof:* Clearly,  $A \preceq 2^A$ , by taking  $f : A \rightarrow 2^A$  to be  $f(x) = \{x\}$ .

To show that  $A \not\approx 2^A$ , we argue by contradiction. Suppose that there *were* a bijection  $f : A \rightarrow 2^A$ . I'll show that this assumption leads to an absurdity.

Construct the following set:

$$A_0 = \{x \in A \mid x \notin f(x)\}$$

This a well-defined subset of  $A$ . Therefore, because  $f$  is onto, there must exist  $a_0 \in A_0$  such that  $f(a_0) = A_0$ .

Now, does  $a_0 \in A_0$ ? There are only two possibilities, yes or no. Neither works:

- If  $a_0 \in A_0$ , then by definition of  $a_0$ ,  $a_0 \notin f(a_0)$ , that is,  $a_0 \notin f(a_0) = A_0$ .
- If  $a_0 \notin A_0$ , then by definition of  $a_0$ ,  $a_0 \in f(a_0)$  (otherwise,  $a_0$  would be in  $A_0$ ) and thus  $a_0 \in f(a_0) = A_0$ .

Either way, we get an absurdity. So our assumption that there is a bijection  $f$  cannot be. Thus,  $A \not\approx 2^A$ . □

Cantor's Theorem means, in particular, that we get an infinite tower of sets of increasing infinite sizes:

$$\mathbb{N} \prec 2^{\mathbb{N}} \prec 2^{2^{\mathbb{N}}} \prec 2^{2^{2^{\mathbb{N}}}} \prec \dots$$

We can now show that there must be an undecidable language. Let  $T(\Sigma)$  be the set of Turing machines over alphabet  $\Sigma$ . We establish that

$$T(\Sigma) \prec 2^{\Sigma^*}$$

where of course  $2^{\Sigma^*}$  is the set of all languages over alphabet  $\Sigma$ .

First, we show that  $T(\Sigma) \preceq \mathbb{N}$ . For this, it suffices to encode every Turing machines  $M$  into a natural number. We identify Turing machines that differ only by the names of the states.

Intuitively, if  $M = (Q, \Sigma, \Gamma, \sqsubset, \vdash, \delta, s, acc, rej)$ , without loss of generality, we can take  $Q = \{1, \dots, n\}$ , and if  $\Gamma = \{a_1, \dots, a_k\}$ , we can encode every symbol as an integer in  $\{1, \dots, k\}$ , taking  $1, \dots, |\Sigma|$  to be the symbols from  $\Sigma$  under a fixed by arbitrary order, and taking  $\sqsubset = |\Sigma| + 1$  and  $\vdash = |\Sigma| + 2$ . This means that we can represent the non-transition parts of  $M$  using numbers  $n, k, s, acc$ , and  $rej$ . For the transition function, note that  $\delta$  can be described by the list of every tuple  $(p, a, q, b, d)$  describing the transition relation, where each tuple is made up of natural numbers, and with  $nk$  tuples total (one for each choice of state and symbol). Thus, we can represent the transition function using  $5nk$  natural numbers, and therefore we can represent  $M$  with natural numbers  $n, k, s, acc, rej$ , and the  $5nk$  natural numbers describing the transition function. We can take these  $5nk + 5$  natural numbers and encode them uniquely in a single natural number using a Goedel encoding: if we let  $2, 3, 5, 7, 11, 13, 17, 19, \dots$  be an enumeration of distinct primes, then we can represent a tuple  $(n_1, n_2, n_3, \dots, n_l)$  as  $2^{n_1} 3^{n_2} 5^{n_3} 7^{n_4} \dots$  (Unique factorization of natural numbers gives us that different Turing machines yield different natural numbers.)

This means that to every Turing machine  $M$  we can associate a distinct natural number, giving us a one-to-one map from  $T(\Sigma)$  to  $\mathbb{N}$ . So  $T(\Sigma) \preceq \mathbb{N}$ .

Next, we can show that  $\mathbb{N} \preceq \Sigma^*$ . That's actually pretty easy. Since  $\Sigma \neq \emptyset$ , then there is an  $a \in \Sigma$ . The map  $f : \mathbb{N} \rightarrow \Sigma^*$  given by  $f(n) = a^n$  is clearly one-to-one. So  $\mathbb{N} \preceq \Sigma^*$ .

Finally, Cantor's Theorem gives us that  $\Sigma^* \prec 2^{\Sigma^*}$ . So we have

$$T(\Sigma) \preceq \mathbb{N} \preceq \Sigma^* \prec 2^{\Sigma^*}$$

and transitivity gives us  $T(\Sigma) \prec 2^{\Sigma^*}$ .

By property (iii), this means that there is no onto function  $T(\Sigma) \rightarrow 2^{\Sigma^*}$ . Consider the function  $L : T(\Sigma) \rightarrow 2^{\Sigma^*}$  that associates to every Turing machine the language it accepts — that  $L$  cannot be onto. Therefore, there must be a language  $A \in 2^{\Sigma^*}$  such that there is no  $M$  with  $L(M) = A$ . That is,  $A$  is undecidable.

## The Halting Problem

The argument above shows that there must be at least one undecidable language. It doesn't help us identify one.

To do so, first we need one of Turing's main result about his machines: universality. It is possible to develop a *universal* Turing machine  $U$  that takes as input (an encoding of) a Turing machine  $M$  and an input string  $w$  and simulates running Turing machine  $M$  on input string  $w$ , accepting when  $M$  accepts, rejecting when  $M$  rejects, and looping when  $M$  loops. I'm not going to give the description of such a Turing machine  $U$ , but I'll put some pointers

on the course web page. Note that this is no stranger than writing, say, a Python interpreter in Python. The key here is that one Turing machine can simulate other Turing machines.

To do so,  $U$  must take a Turing machine as input, that is, we need a way to encode Turing machines as input of a given alphabet. For the sake of concreteness, let's fix  $\Sigma = \{0, 1\}$  here. (The undecidability result does not depend on the alphabet.) It is clear to our computer scientist intuition that we can represent a Turing machine as a string of bits. We need to be able to read off the various components of an encoded Turing machine, but we can do so by encoding each component separately. We've done so with natural numbers above. It's a simple matter to do so with strings of bits. The details of the encoding are unimportant, so I'll just assume that if  $M$  is a Turing machine, there is an encoding  $\widehat{M}$  of that Turing machine. Similarly, let  $\langle u, v \rangle$  be an encoding of the pair  $u$  and  $v$  in such a way that we can recover both  $u$  and  $v$  from the encoding.

Thus, universal Turing machine  $U$  takes as input strings of the form  $\langle \widehat{M}, w \rangle$  — encodings of a pair of a Turing machine encoding and a string, and simulates  $M$  with input  $w$ .

Consider the following language, called the *Halting Problem*:

$$HP = \{ \langle \widehat{M}, w \rangle \mid M \text{ halts on input } w \}$$

I claim that  $HP$  is undecidable, that is, there is no total Turing machine  $M$  that accepts  $HP$ . It turns out to be a similar argument than the one used for Cantor's Theorem. We argue by contradiction: assume it *is* decidable, and derive an absurdity.

By way of contradiction, assume  $HP$  is decidable. That means we have a total Turing machine  $K$  that accepts  $HP$ . That is,  $K$  accepts  $\langle \widehat{M}, w \rangle$  when  $M$  halts on  $w$ , and rejects when  $M$  does not halt on  $w$ .

Using  $K$ , construct another Turing machine  $I$  as follows:

On input  $x$ :

1. Run  $U$  with input  $\langle \widehat{K}, \langle x, x \rangle \rangle$
2. If  $U$  rejects, accept
3. If  $U$  accepts, go into an infinite loop

Clearly, we can implement  $I$  as a Turing machine if we have  $K$  and  $U$ . Since  $I$  is a Turing machine, it has an encoding  $\widehat{I}$ . We can also ask whether  $I$  halts on a given input. The input we care about?  $\widehat{I}$  of course!

There are only two possibilities. Either  $I$  halts on input  $\widehat{I}$ , or it does not. Neither makes sense.

- Say  $I$  halts on input  $\widehat{I}$ . By definition of  $I$ , this happens only when  $U$  rejects  $\langle \widehat{K}, \langle \widehat{I}, \widehat{I} \rangle \rangle$ . Since  $U$  is a universal Turing machine,  $U$  rejects when  $K$  rejects  $\langle \widehat{I}, \widehat{I} \rangle$ . But  $K$  is the Turing machine deciding  $HP$ , and it rejects exactly when  $I$  does not halt on  $\widehat{I}$ . But we said  $I$  halts on  $\widehat{I}$ . That's absurd.
- Say  $I$  does not halt on input  $\widehat{I}$ . By definition of  $I$ , this happens only when  $U$  accepts

$\langle \widehat{K}, \langle \widehat{I}, \widehat{I} \rangle \rangle$ . But  $U$  is a universal Turing machine, so  $U$  accepts when  $K$  accepts  $\langle \widehat{I}, \widehat{I} \rangle$ . But  $K$  is the Turing machine deciding  $HP$ , and it accepts when  $I$  halts on input  $\widehat{I}$ . But we said  $I$  does not halt on  $\widehat{I}$ . That's absurd.

Either way, we get an absurdity. So our assumption that  $HP$  is decidable, that is, that  $K$  exists, must be wrong. There is no such  $K$ . So  $HP$  is undecidable.