

## Notes on Distinguishability

This notes present a technique to prove a lower bound on the number of states of any DFA that recognizes a given language. The technique can also be used to prove that a language is not regular. (By showing that *for every*  $k$  one needs at least  $k$  states to recognize the language.)

It will be helpful to keep in mind the following two languages over the alphabet  $\Sigma = \{0, 1\}$  as examples: the language  $EQ = \{0^n 1^n | n \geq 1\}$  of strings containing a sequence of zeroes followed by an equally long sequence of ones, and the language  $A = (0 \cup 1)^* \cdot 1 \cdot (0 \cup 1)$  of strings containing a 1 in the second-to-last position.

We start with the following basic notion.

**Definition 1 (Distinguishable Strings)** *Let  $L$  be a language over an alphabet  $\Sigma$ . We say that two strings  $x$  and  $y$  are **distinguishable** with respect to  $L$  if there is a string  $z$  such that  $xz \in L$  and  $yz \notin L$ , or vice versa.*

For example the strings  $x = 0$  and  $y = 00$  are distinguishable with respect to  $EQ$ , because if we take  $z = 1$  we have  $xz = 01 \in EQ$  and  $yz = 001 \notin L$ . Also, the strings  $x = 00$  and  $y = 01$  are distinguishable with respect to  $A$  as can be seen by taking  $z = 0$ .

On the other hand, the strings  $x = 0110$  and  $y = 10$  are *not* distinguishable with respect to  $EQ$  because for every  $z$  we have  $xz \notin L$  and  $yz \notin L$ .

**Exercise 1** *Find two strings that are not distinguishable with respect to  $A$ .*

The intuition behind Definition 1 is captured by the following simple fact.

**Lemma 1** *Let  $L$  be a language,  $M$  be a DFA that decides  $L$ , and  $x$  and  $y$  be distinguishable strings with respect to  $L$ . Then the state reached by  $M$  on input  $x$  is different from the state reached by  $M$  on input  $y$ .*

PROOF: Suppose by contradiction that  $M$  reaches the same state  $q$  on input  $x$  and on input  $y$ . Let  $z$  be the string such that  $xz \in L$  and  $yz \notin L$  (or vice versa). Let us call  $q'$  the state reached by  $M$  on input  $xz$ . Note that  $q'$  is the state reached by  $M$  starting from  $q$  and given the string  $z$ . But also, on input  $yz$ ,  $M$  must reach the same state  $q'$ , because  $M$  reaches state  $q$  given  $y$ , and then goes from  $q$  to  $q'$  given  $z$ . This means that  $M$  either accepts both  $xz$  and  $yz$ , or it rejects both. In either case,  $M$  is incorrect and we reach a contradiction.  $\square$

Consider now the following generalization of the notion of distinguishability.

**Definition 2 (Distinguishable Set of Strings)** *Let  $L$  be a language. A set of strings  $\{x_1, \dots, x_k\}$  is distinguishable if for every two distinct strings  $x_i, x_j$  we have that  $x_i$  is distinguishable from  $x_j$ .*

For example one can verify that  $\{0, 00, 000\}$  are distinguishable with respect to  $EQ$  and that  $\{00, 01, 10, 11\}$  are distinguishable with respect to  $A$ .

We now prove the main result of this handout.

**Lemma 2 (Main)** *Let  $L$  be a language, and suppose there is a set of  $k$  distinguishable strings with respect to  $L$ . Then every DFA for  $L$  has at least  $k$  states.*

PROOF: If  $L$  is not regular, then there is no DFA for  $L$ , much less a DFA with less than  $k$  states. If  $L$  is regular, let  $M$  be a DFA for  $L$ , let  $x_1, \dots, x_k$  be the distinguishable strings, and let  $q_i$  be the state reached by  $M$  after reading  $x_i$ . For every  $i \neq j$ , we have that  $x_i$  and  $x_j$  are distinguishable, and so  $q_i \neq q_j$  because of Lemma 1. So we have  $k$  different states  $q_1, \dots, q_k$  in  $M$ , and so  $M$  has at least  $k$  states.  $\square$

Using Lemma 2 and the fact that the strings  $\{00, 01, 10, 11\}$  are distinguishable with respect to  $A$  we conclude that every DFA for  $A$  has at least 4 states.

For every  $k \geq 1$ , consider the set  $\{0, 00, \dots, 0^k\}$  of strings made of  $k$  or fewer zeroes. It is easy to see that this is a set of distinguishable strings with respect to  $EQ$ . This means that there cannot be a DFA for  $EQ$ , because, if there were one, it would have to have at least  $k$  states for every  $k$ , which is clearly impossible.