

CMPSCI 250: Introduction to Computation

Lecture #31: What DFA's Can and Can't Do
David Mix Barrington
9 April 2014

What DFA's Can and Can't Do

- Deterministic Finite Automata
- Formal Definition of DFA's
- Examples of DFA's
- DFA's in Java
- Characterizing Strings With Given Behavior
- Distinguishable Strings
- Languages With No DFA's

Deterministic Finite Automata

- We now turn to **finite-state machines**, a model of computation that captures the idea of reading a file of text with a fixed limit on the memory we can use to remember what we have seen.
- In particular, the memory used must be **constant**, independent of the length of the file. We ensure this by requiring our machine to have a **finite state set**, so that at any time during the computation all that it knows is which state it is in.

Deterministic Finite Automata

- The **initial state** is fixed. When the machine sees a new letter, it changes to a new state based on a fixed **transition function**. When it finishes the string, it gives a yes or no answer based on whether it is in a **final state**.
- Because the new state depends only on the old state and the letter seen, the computation is **deterministic** and the machine is called a **deterministic finite automaton** or **DFA**.

Where We Are Going

- A DFA **decides** a language -- it reads a string over its alphabet and then answers “yes” or “no”. The **language of the DFA** is the set of strings for which it says yes.
- We call a language X **decidable** if there exists a DFA whose language is X . Later we’ll prove that some languages are *not* decidable.

Where We Are Going

- The **Myhill-Nerode Theorem** will give us a way to take an arbitrary language and determine whether it is decidable.
- We'll define a particular equivalence relation on strings, based only on the language. If this relation has a finite set of equivalence classes, there is a DFA for the language, and there is a **minimal DFA** with as many states as there are classes. We'll see how to compute the minimal DFA from any DFA for the language.

Where We Are Going

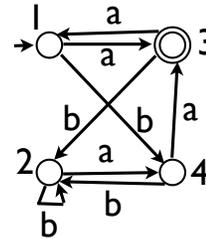
- As we've mentioned, there is a DFA for a language if and only if the language is **regular** (that is, if and only if it is the language denoted by some regular expression).
- We'll prove this important result, called **Kleene's Theorem**, over several lectures. Our proofs will show us how to convert a DFA to a regular expression and vice versa.

Formal Definition of DFA's

- Formally a DFA is defined by its **state set** S , its **initial state** $i \in S$, its **final state set** $F \subseteq S$, its **input alphabet** Σ , and its **transition function** δ from $(S \times \Sigma)$ to S .
- We usually represent DFA's by diagrams (labeled directed multigraphs) with a node for each state, a special mark for the initial state, a double circle on each final state, and an arrow labeled "a" from node p to node q whenever $\delta(p, a) = q$.

A DFA Example

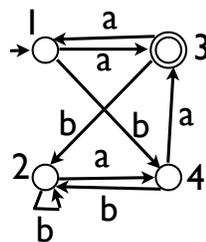
- Here is a DFA with state set $\{1, 2, 3, 4\}$, initial state 1, final state set $\{3\}$, alphabet $\{a, b\}$, and a transition function indicated by the arrows.



- For any string, we can follow the arrows for its letters in order. The strings a , ba , and $bbbaa$ are in this DFA's language.

Clicker Question #1

- Which of these strings is not in the language of this DFA?
- (a) abababbaaab
- (b) bababaaaba
- (c) aaaabbbbbaaaa
- (d) bbbbbabaabaa



Answer #1

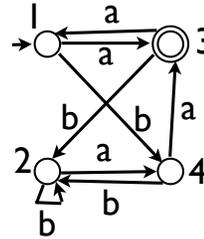
- Which of these strings is not in the language of this DFA?

- (a) *abababbaaab*
(1-3-2-4-2-4-2-2-4-3-1-4)

- (b) bababaaaba
(1-4-3-2-4-2-4-3-1-4-3)

- (c) aaaabbbbbaaaa
(1-3-1-3-1-4-2-2-2-4-3-1-3)

- (d) bbbbbabaabaa
(1-4-2-2-2-2-4-2-4-3-2-4-3)

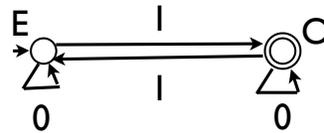


Behavior of a DFA

- The **behavior function** of a particular DFA is a function called δ^* from $(S \times \Sigma^*)$ to S , such that $\delta^*(p, w)$ is the state of the DFA after it starts in state p and reads the string w .
- Formally, we say that $\delta^*(p, \lambda) = p$ and that $\delta^*(p, wa) = \delta(\delta^*(p, w), a)$.
- The **language of a DFA** is defined to be the set of strings w such that $\delta^*(i, w)$ is a final state. For a DFA M , we call this language $L(M)$.

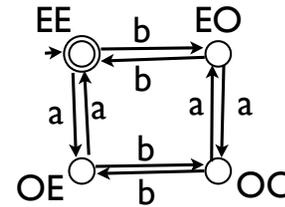
More Examples of DFA's

- One of the simplest possible DFA's decides the language of binary strings with an odd number of ones. It has two states E and O, representing whether the machine has seen an even or odd number of ones so far. The initial state is E, and the final state set is {O}. The transition function has $\delta(E, 0) = E$, $\delta(E, 1) = O$, $\delta(O, 0) = O$, and $\delta(O, 1) = E$.



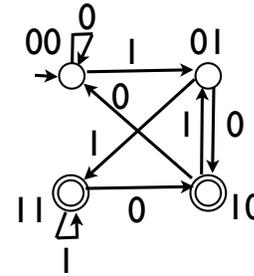
More Examples of DFA's

- We can build a four-state DFA for the language EE , of strings with an even number of a's and an even number of b's.
- Its states are EE , EO , OE , and OO . For example, $\delta^*(EE, w) = EO$ if w has an even number of a's and an odd number of b's.



More Examples of DFA's

- Another four-state DFA can decide whether the next to last letter of a binary string w exists and is 1.
- The state set is $\{00, 01, 10, 11\}$ and the state after reading w represents the last two letters seen. The initial state is 00 and the final state set is $\{10, 11\}$.



DFA's in Pseudo-Java

- We consider the input to be given like a file, with a method to give the next letter and one to tell when the input is done.
- We relabel the state set and the alphabet to be $\{0, \dots, \text{states} - 1\}$ and $\{0, \dots, \text{letters} - 1\}$ respectively.

DFA's in Pseudo-Java

```
public class DFA {
    natural states;
    natural letters;
    natural start;
    boolean [ ] isFinal =
        new boolean[states];
    natural [ ] [ ] delta =
        new natural [states] [letters];
    natural getNext( ) {code omitted}
    boolean inputDone( ) {code omitted}
```

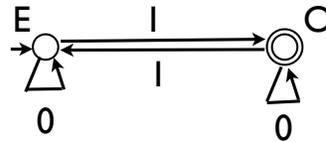
DFA's in Pseudo-Java

```
boolean decide ( )
{ // returns whether input string is
  // in the language of the DFA

  natural current = start;
  while (!inputDone( ))
    current =
      delta[current][getNext( )];
  return isFinal [current];}
```

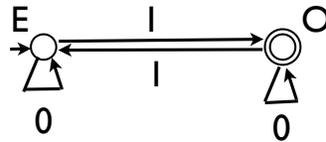
The Strings With a Behavior

- How do we prove that a particular DFA has a particular language?
- With the even-odd DFA, we can say that $\delta^*(E, w) = E$ if w has an even number of ones, and $\delta^*(E, w) = O$ if it has an odd number of ones.



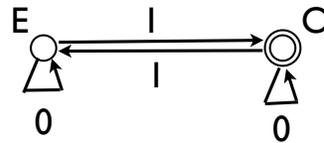
The Strings With a Behavior

- “ $\delta^*(E, w) = E$ if w has an even number of ones, and $\delta^*(E, w) = O$ if it has an odd number of ones.”
- Letting $P(w)$ be the entire statement in the bullet above, we can prove $\forall w:P(w)$ by induction on all binary strings. $P(\lambda)$ says that $\delta^*(E, \lambda) = E$, because λ has no ones and 0 is even, and $\delta^*(E, \lambda) = E$ is true by definition of δ^* .



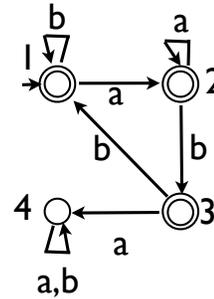
The Strings With a Behavior

- Now assume that $P(w)$ is true, so that $\delta^*(E, w)$ is E if w has an even number of ones and O otherwise. Then $w0$ has the same number of ones as w , so $\delta^*(E, w0)$ should be the same state as $\delta^*(E, w)$. And $w1$ has one more one than w , so $\delta^*(E, w1)$ should be the other state from $\delta^*(E, w)$. In each of the four cases, the new state is the state given by the δ function of the DFA.



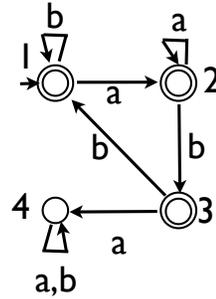
The No-aba Language

- The language No-aba is the set of strings that never have an aba substring.
- We can build a DFA M for No-aba with state set $\{1, 2, 3, 4\}$, start state 1, final state set $\{1, 2, 3\}$, and transition function as shown. (We call 4 a **death state**.) We can see that an aba will take us from any state to 4.



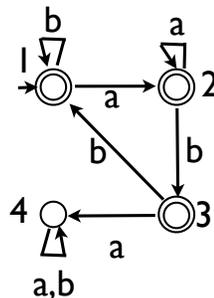
Clicker Question #2

- Suppose $\delta^*(1, w) = 4$. Which statement must be true of w ?
- (a) w contains an aba substring
- (b) $w \in \text{No-aba}$
- (c) w contains no aaa
- (d) $w \in (\text{abb})^* \text{aba}(\text{a+b})^*$



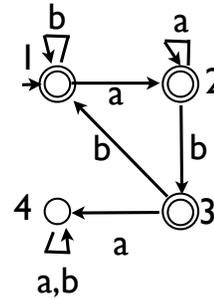
Answer #2

- Suppose $\delta^*(1, w) = 4$.
Which statement must be true of w ?
- (a) w contains an *aba* substring
- (b) $w \in \text{No-aba}$
- (c) w contains no *aaa*
- (d) $w \in (\text{abb})^* \text{aba}(\text{a+b})^*$



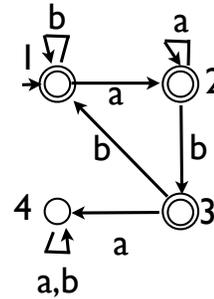
Characterizing the States

- Let L_1 be the set of strings that have no aba and don't end in a or ab.
- Let L_2 be the set of strings that don't have an aba and end in a.
- L_3 is the set of strings that don't have an aba and end in ab.
- L_4 is the set that have aba.



Characterizing the States

- We can make eight checks, one for each value of δ . If $\delta(i, x) = j$, we check that any string in L_i , followed by the letter x , yields a string in L_j .
- We then do an inductive proof, where $P(w)$ is the statement on the previous slide: “For all states i , $\delta^*(i, w) = i$ if and only if $w \in L_i$ ” where each L_i is as defined there.
Thus $w \in L(M) \leftrightarrow w \in \text{No-aba}$.



Distinguishable Strings

- Is it possible that another DFA with only three states could decide No-aba?
- We divided all possible strings into four sets. Suppose a DFA reads w and does not know which of the four sets w is in.
- We'll show that in this case it is *doomed* -- for some string x , it will be wrong if it sees x and has to decide whether wx is in the language No-aba.

Distinguishing Strings

- Look at the four strings λ , a , ab , and aba .
- If the DFA has $\delta^*(i, \lambda) = \delta^*(i, a)$, we say that it **cannot distinguish between** λ and a .
- If this is true, the DFA must also have $\delta^*(i, b) = \delta^*(i, ab)$ because a b will take the same state to the same state.
- Then as well $\delta^*(i, ba) = \delta^*(i, aba)$.

Distinguishable Strings

- But now we know that the DFA cannot decide No-aba, because it gives the same answer on the strings ba (which is in No-aba) and aba (which is not in No-aba).
- We can call this an **experiment** that **distinguishes** the two strings λ and a.

Clicker Question #3

- Two strings u and v are defined to be No-aba-distinguishable if there exists a string w such that exactly one of the strings uw and vw are in No-aba. Which one of these pairs of strings is not No-aba-distinguishable?
- (a) {aabab, bababa}
- (b) $\{\lambda, babbbaab\}$
- (c) {abbabba, abb}
- (d) {abbabbaaba, babbaa}

Answer #3

- Two strings u and v are defined to be No-aba-distinguishable if there exists a string w such that exactly one of the strings uw and vw are in No-aba. Which one of these pairs of strings is No-aba-distinguishable?
- (a) $\{aabab, bababa\}$ (both already have aba)
- (b) $\{\lambda, babbbaab\}$ (append a to each)
- (c) $\{abbabba, abb\}$ (append ba to each)
- (d) $\{abbabbaaba, babbaa\}$ (append λ to each)

Sets of Distinguishable Strings

- Let L be any language. We say that two strings u and v are **L-distinguishable** (also called **L-inequivalent**) if there exists a string w such that $uw \in L$ and $vw \notin L$, or vice versa.
- We call the strings **L-equivalent** if the negation of this statement is true, that is, if $\forall w: uw \in L \leftrightarrow vw \in L$.

A Lemma on Distinguishability

- **Lemma:** If M is a DFA with transition function δ , L is any language, u and v are two L -distinguishable strings, and $\delta^*(i, u) = \delta^*(i, v)$, then $L(M) \neq L$.
- **Proof:** We can prove by induction that if $\delta^*(i, u) = \delta^*(i, v)$, then for any string w , $\delta^*(i, uw) = \delta^*(i, vw)$. For the particular w that distinguishes u and v , then, the single state $\delta^*(i, uw) = \delta^*(i, vw)$ needs to be both final and non-final if $L(M) = L$.

A Distinguishability Theorem

- **Theorem:** If there exists a set of k *pairwise L-distinguishable* strings, then no DFA that decides L can have fewer than k states.
- **Proof:** If there are more strings than there are states, by the **Pigeonhole Principle** there must exist two L -distinguishable strings u and v such that $\delta^*(i, u) = \delta^*(i, v)$. In this case the Lemma says that the DFA does not decide L .

Languages With No DFA's

- Consider the balanced parenthesis language Paren, which we will write as a subset of $\{L, R\}^*$ with L for left parens and R for right parens. We can prove that there is *no DFA at all* that decides this language.
- Look at the infinite set of strings $\{\lambda, L, LL, LLL, \dots\}$. I claim that this set is pairwise Paren-distinguishable, because if i and j are two naturals with $i \neq j$, then L^i and L^j are distinguished by R^i , since $L^i R^i$ is in Paren and $L^j R^i$ is not.

Languages With No DFA's

- So for any natural k , we can find more than k pairwise Paren-distinguishable strings, and by our theorem there cannot be a k -state DFA.
- Our real-life algorithm to decide Paren is to remember the number of L's we have seen, minus the number of R's. If this number ends at 0, without ever going negative, we are in Paren. But this requires more than constant memory -- potentially a state for every natural from 0 through n .