

Regular Languages and Automata: Part 1

1 Regular expressions

See also the video “Regular Expressions” on Canvas (16 minutes).

Look at these two problems.

1. A string is a “valid password” when it contains at least 8 characters and at least 2 digits. Given a string, say whether it is a valid password.
2. Given a string (e.g. a file), list all occurrences of email addresses within it. Each occurrence should be represented as a pair of numbers (i, m) , where i is the start position (e.g. 0 for an occurrence at the start of the file) and m is the length.

The first of these is an example of a *matching problem*. The second is a *finding problem*. Such problems—and variations—often arise in computing, so people have made tools to solve them efficiently. To use such a tool, you have to specify when a word is a valid password or an email address or whatever. Often, we do this by means of a *regular expression*.

1.1 Definitions

The alphabet (set of characters) is called Σ . To keep things simple, let’s suppose that it’s $\{a, b, c\}$. In a practical situation, it might instead be the ASCII alphabet, which has 128 characters. Or it might be the Unicode alphabet, which has 137,439 characters. In any case, we’ll assume that Σ is a finite set and contains at least two characters.

We write Σ^* for the set of all words. A *language* is a set of words, i.e. a subset of Σ^* . For example: the set of valid passwords is a language, and so is the set of email addresses. A regular expression (regex), such as $c(bb|ca)^*$, represents a language, just as an arithmetic expression, such as $2 + (5 \times 3)$, represents a number.

Now let me explain how regexps work.

- The regexp a matches only the word a .
- The regexp b matches only the word b .
- The regexp c matches only the word c .
- The regexp ε matches only the empty word ε .
- If E and F are regexps, then the regexp EF matches any word that’s a concatenation of a word matched by E and a word matched by F .
- If E and F are regexps, then the regexp $E|F$ matches any word that either E or F matches.
- If E is a regexp, then the regexp E^* matches any word that’s a concatenation of several (i.e. zero or more) words matched by E .
- The (rarely used) regexp 0 doesn’t match any word.

1.2 Precedence

For arithmetic expressions, \times has higher precedence than $+$, and knowing this enables us to parse the expression $3 + 4 \times 2$ as $3 + (4 \times 2)$, for example. For regexps, the precedence laws are as follows: juxtaposition (which means “putting things next to each other”) has higher precedence than $|$ and lower precedence than $*$. Knowing this enables us to parse $c(bb|ca)^*$ as $c(((bb)|(ca))^*)$, for example.

More operators that you can use:

- E^+ is short for EE^* . It matches any word that is a concatenation of one or more words matched by E .
- $E?$ is short for $\varepsilon|E$.

These have the same precedence as $*$.

Some tools provide additional operators, which make it possible to express fancier languages. Expressions using these additional operators may be called “regular expressions” in the tool documentation, but technically they are not regular.

- Example 1** 1. Does the regexp $c(bb|ca)^*$ match $ccacabb$? YES/NO.
2. Does the regexp $c(bb|ca)^*$ match $cbbcacac$? YES/NO.
3. Does the regexp $(c(bb|ca)^*)^*$ match $cccacacbbcbba$? YES/NO.
4. Do $(a|b)c^*$ and $ac^*|bc^*$ represent the same language? YES/NO.
5. Do $(a|b)c^*$ and ac^* represent the same language? YES/NO.

2 Some questions about regular expressions

Before we look into regexps in more detail, let’s consider some questions. For some of these, the answers are far from obvious, but will emerge over the coming lectures.

2.1 Regular and Irregular Languages

Recall that a regexp represents a language. Any language $L \subseteq \Sigma^*$ that can be represented in this way is said to be *regular*. Questions:

1. Are there any languages that are not regular? Answer: Yes, and we’ll see some examples.
2. Is the complement of a regular language always regular? (For example, is there a regexp for those words that are *not* matched by $c(bb|ca)^*$?) Answer: Yes.
3. Is the intersection of two regular languages always regular? (For example, is there a regexp for those words that are matched by *both* $cc(bb|ca)^*$ and $c(bbbb|cca)^*$?) Answer: Yes.

2.2 Decidability questions

A *decision problem* is a problem that, for any given argument, has a Yes/No answer. For example, the finding problem above is not a decision problem, because the answer (for a given file) is a set of pairs of numbers. A decision problem is said to be *decidable* when there is some program that, given an argument, says whether the answer is Yes or No.

We can ask the following questions about regexps:

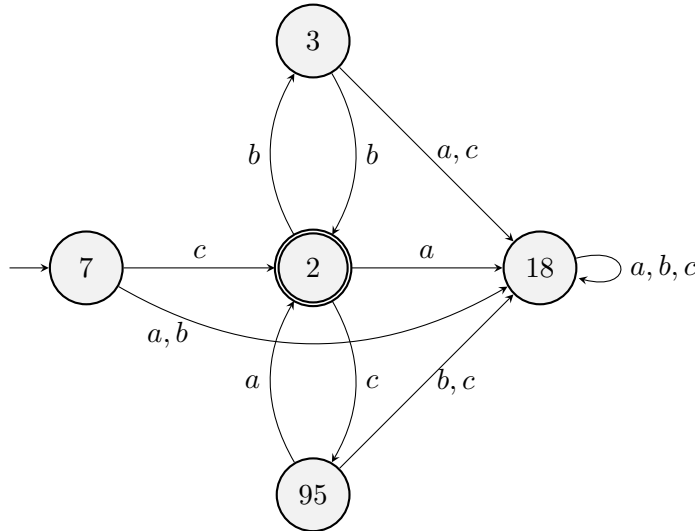
1. Is the matching problem for the regexp $c(bb|ca)^*$ *decidable*? In other words, is there some program that, when given a word w over our alphabet $\Sigma = \{a, b, c\}$, returns True if w matches $c(bb|ca)^*$, and False if it doesn’t? (If w isn’t a word over our alphabet, then it doesn’t matter what happens.) Answer: Yes.
2. Is the matching problem for the regexp $(c(bb|ca)^*)^*$ decidable? Answer: Yes.
3. Is it the case that, for *every* regexp E , the matching problem for E is decidable? Answer: Yes.
4. Is the matching problem for regexps decidable? In other words, is there some program that, when given a regexp E and word w , returns True if w matches E , and False if it doesn’t? Answer: Yes.
5. Is language equality for regexps decidable? In other words, is there some program that, when given regexps E and F , returns True if they represent the same language and False otherwise? Answer: Yes.

2.3 Efficiency questions

In the previous section, we asked whether certain problems can be solved at all. Another question is: can they be solved efficiently? After all, your customers aren't willing to wait a long time for an answer from your program. This question isn't very precise, but it's important. We'll see that for some of these problems, we can give a reasonably efficient solution.

3 Introducing automata

See also the video "Introduction to automata" on Canvas (10 minutes).



Recall that we wanted a program to solve the matching problem for $c(bb|ca)^*$. This can be achieved by the automaton shown. There are five *states*, represented as circles. The automaton processes a word by starting at the *initial state* (indicated by \rightarrow) and performing a *transition* as it inputs each letter. When the whole word has been input, the automaton returns Yes if the current state is *accepting*, indicated by a double ring; otherwise it returns No.

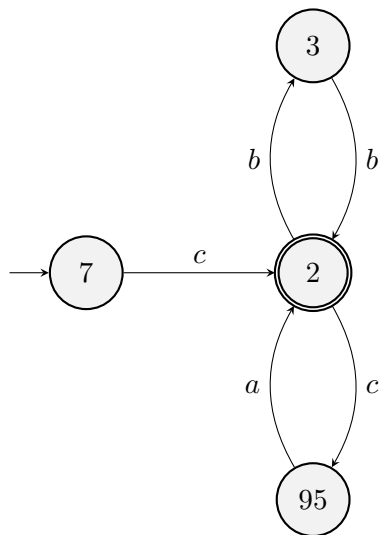
This is a *deterministic finite automaton* (DFA). "Deterministic" because the initial state and the result of each transition are specified. "Finite" because the set of states is finite.

Definition 1 A deterministic finite automaton consists of the following data.

- A finite set X of states.
- An initial state $p \in X$.
- A transition function $\delta: X \times \Sigma \rightarrow X$.
- A set of accepting states $\text{Acc} \subseteq X$.

In the above example,

$$\begin{aligned}
 X &= \{7, 2, 3, 95, 18\} \\
 p &= 7 \\
 \delta &= \{(7, a) \mapsto 18, (7, b) \mapsto 18, (7, c) \mapsto 2, \\
 &\quad (2, a) \mapsto 18, (2, b) \mapsto 3, (2, c) \mapsto 95, \\
 &\quad (3, a) \mapsto 18, (3, b) \mapsto 2, (3, c) \mapsto 18, \\
 &\quad (95, a) \mapsto 2, (95, b) \mapsto 18, (95, c) \mapsto 18, \\
 &\quad (18, a) \mapsto 18, (18, b) \mapsto 18, (18, c) \mapsto 18\} \\
 \text{Acc} &= \{2\}
 \end{aligned}$$

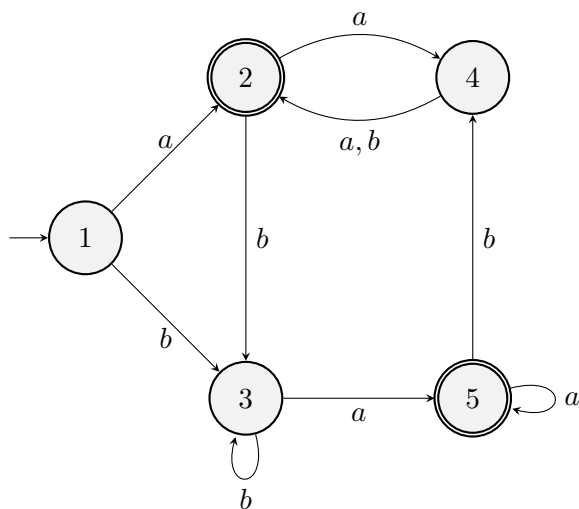


This automaton is more efficient than the one above. It is a *partial DFA*, meaning that δ is merely a partial function, i.e. it can sometimes be undefined. As soon as a character cannot be input, the word is rejected. For example, the word `cbccabbababccabcc` is rejected after just three characters. (A partial DFA can also have no initial state, but then every word is rejected straight away, so this isn't very useful.)

We can easily turn a partial DFA into a total DFA; just add an extra non-accepting state, called the *error state* (18 in the example). Transitions that are undefined in the partial DFA go to the error state. And every transition from the error state goes to the error state. (If the partial DFA has no initial state, the error state will be initial.)

In the above examples, the states are numbers, but it doesn't actually matter what they are. So, when drawing an automaton, it's fine to leave the circles blank, unless we want to refer to a specific state.

Example 2 Here is a DFA over the alphabet $\{a, b\}$. Are these words accepted?



abab	Y/N
ababba	Y/N
ababbbaa	Y/N
bab	Y/N
baa	Y/N
ϵ	Y/N

Example 3 What about this DFA?

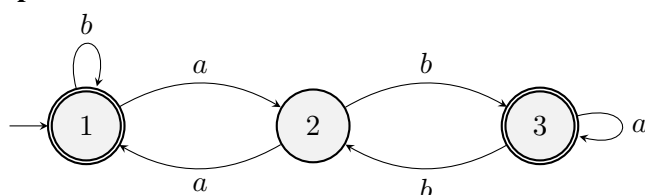
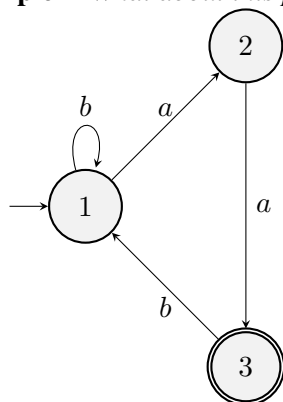


abb	Y/N
abbaba	Y/N
bba	Y/N
ϵ	Y/N

Example 4 What about this partial DFA?

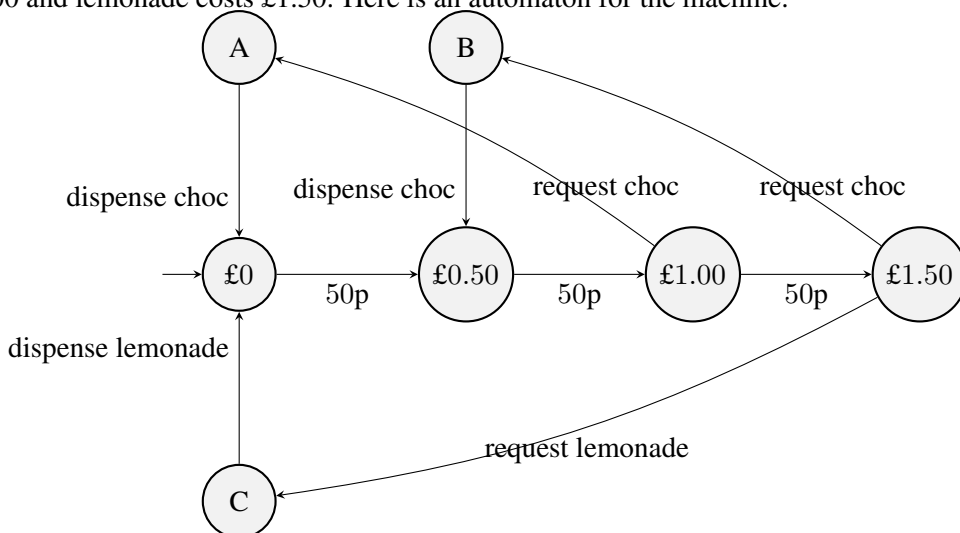


aaab	Y/N
aabaa	Y/N
abaab	Y/N
aa	Y/N

3.1 Vending machines

The idea of a DFA was invented for a specific purpose: solving a language's matching problem. All a DFA can do is input letters and say whether the word that's been read is accepted or not. But for other purposes, there are other kinds of automaton.

In the lobby there's a vending machine that can receive 50p coins, up to a maximum of £1.50. Chocolate costs £1.00 and lemonade costs £1.50. Here is an automaton for the machine:

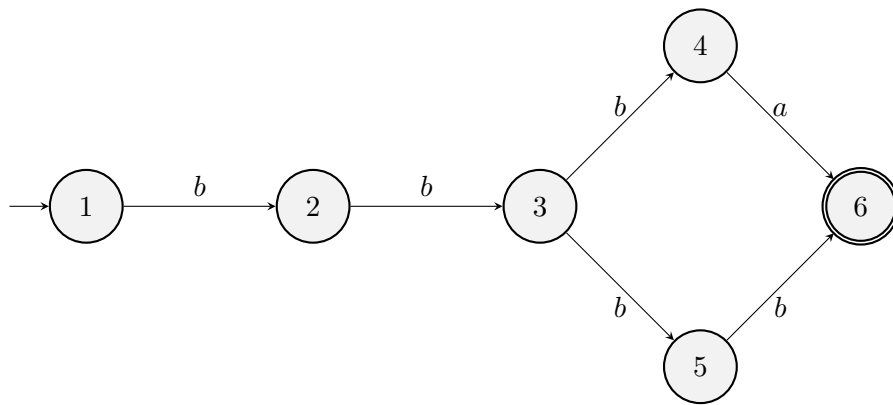


Note that this automaton can *input* money and requests, and also *output* chocolate and lemonade. There are no accepting states, since recognizing words is not the purpose of this machine.

4 Nondeterministic automata

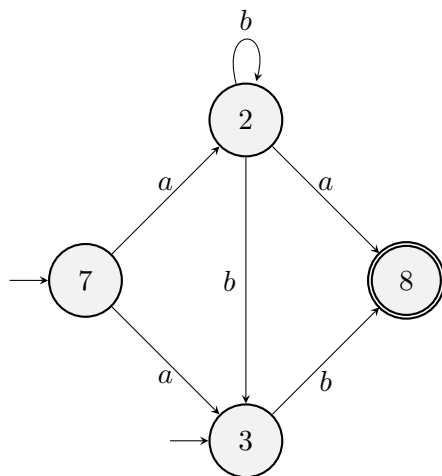
See also the video “Nondeterministic finite automata” on Canvas (6 minutes).

Sometimes it is difficult to obtain a DFA for a regexp, but we can more easily obtain a *nondeterministic* finite automaton (NFA). Here's an example, for the language $bb(ba|bb)$.



A nondeterministic finite automaton (NFA) differs from a DFA in two respects. Firstly an NFA can have several initial states. Secondly, from a given state, when a (or any other character) is input, there can be several possible next states. Thus δ is a *relation* but not a function. The automaton chooses its initial state, and chooses what state to move to as it inputs a character. A word w is *acceptable* when there is **some** path from an initial state to an accepting state that goes through the characters of w .

Example 5 Here's a NFA for the alphabet $\{a, b\}$. Are these words acceptable?

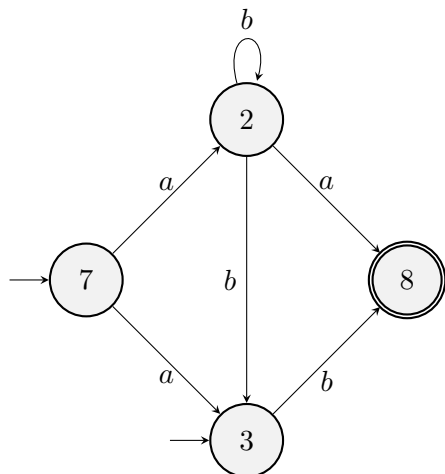


abbb	Y/N
abb	Y/N
ε	Y/N
abba	Y/N

4.1 Determinizing an NFA: transforming an NFA into a DFA

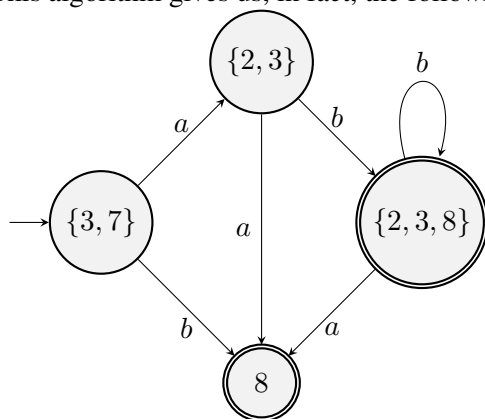
See also the video “Determinizing a NFA” on Canvas (10 minutes).

An NFA is useless in practice: we want a program that always says Yes to a good word and No to a bad word, and an NFA doesn't do that. But we can *determinize* it, i.e. turn it into a DFA that recognizes the same language. To see how this works, look at the following example.



Let's think how to find out whether the word `abb` is acceptable. We can do it by trial and error, but there's an algorithmic way: *keep track of the current set of possible states*. Initially the set of possible states is $\{3, 7\}$. After inputting `a`, the set of possible states is $\{2, 3\}$. After inputting `b`, the set of possible states is $\{2, 3, 8\}$. After inputting `b` again, the set of possible states is $\{2, 3, 8\}$. We have reached the end of the word, and we note that one of the currently possible states, viz. 8, is accepting. Therefore the word `abb` is acceptable.

This algorithm gives us, in fact, the following DFA:



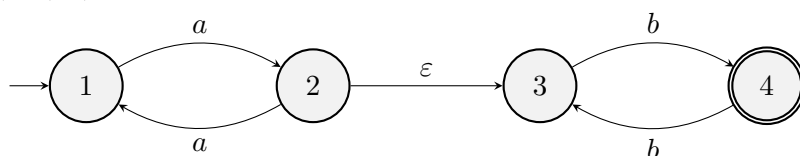
The “states” of the DFA are *sets* of states of the NFA. The initial “state” is the set of all the initial states of the NFA. A “state” is accepting when it contains an accepting state of the NFA. From a given “state”, when we input a character, we collect all the possible next states. This process is called *determinization*.

We see that a word w is accepted by the DFA iff it's acceptable to the NFA. Therefore they represent the same language.

5 ε -transitions

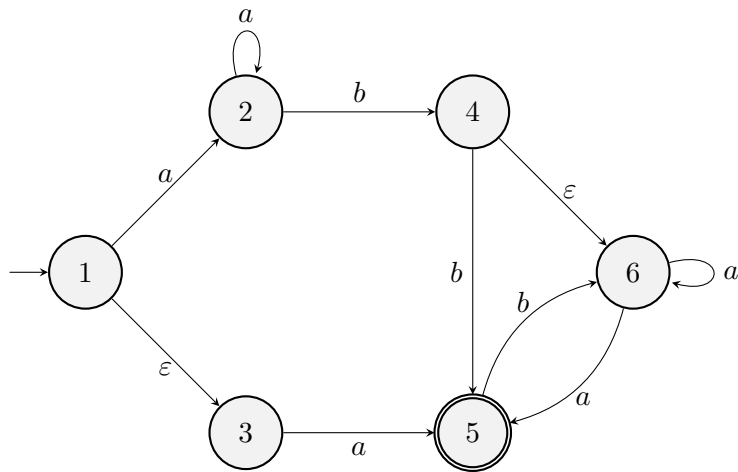
See also the video “Epsilon transitions” on Canvas (6 minutes).

Sometimes even an NFA is difficult to obtain, but we can obtain an automaton that spends some time thinking. As it thinks, it moves from one state to another **without inputting any character**. Here's an example, for the regexp $a(aa)^*b(bb)^*$.



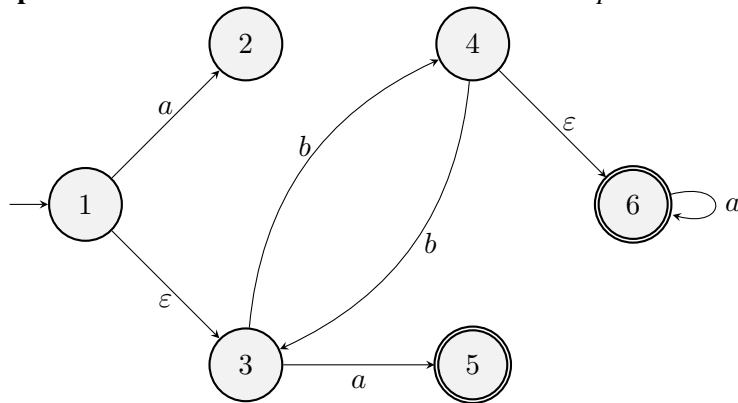
We call this a *nondeterministic automaton with ε -transitions* or ε NFA for short. A word w is acceptable when there is some path from an initial state to an accepting state that goes through the characters of w , padded with ε transitions.

Example 6 Here is a ε NFA. Are these words acceptable?



aba	Y/N
ab	Y/N
aaabb	Y/N
a	Y/N

Example 7 Here's another ϵ NFA. Are these words acceptable?

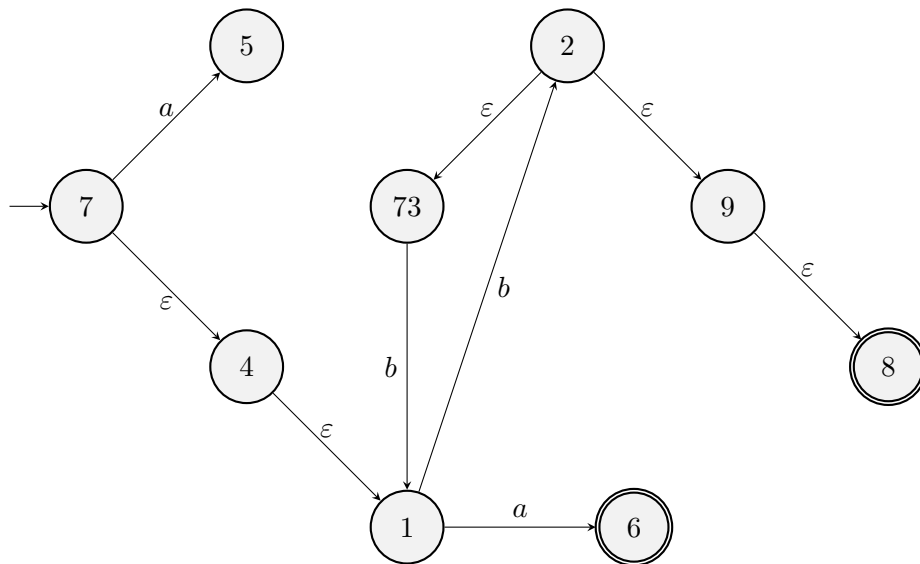


a	Y/N
aba	Y/N
bbb	Y/N
bbba	Y/N
ab	Y/N

5.1 Removing ϵ -transitions

See also the video “Removing epsilon transitions” on Canvas (9 minutes).

Happily, it's possible to remove the ϵ -transitions from a ϵ NFA, i.e. convert it to an NFA that recognizes the same language. To get the idea for this procedure, let's look at the following ϵ NFA.



The word bbb is acceptable, because of the following path:

$7 \xrightarrow{\varepsilon} 4 \xrightarrow{\varepsilon} 1 \xrightarrow{b} 2 \xrightarrow{\varepsilon} 73 \xrightarrow{b} 1 \xrightarrow{b} 2 \xrightarrow{\varepsilon} 9 \xrightarrow{\varepsilon} 8$

This path consists of the following pieces:

$7 \xrightarrow{\varepsilon} 4 \xrightarrow{\varepsilon} 1 \xrightarrow{b} 2$ is a *slow b-transition* from 7 to 2;

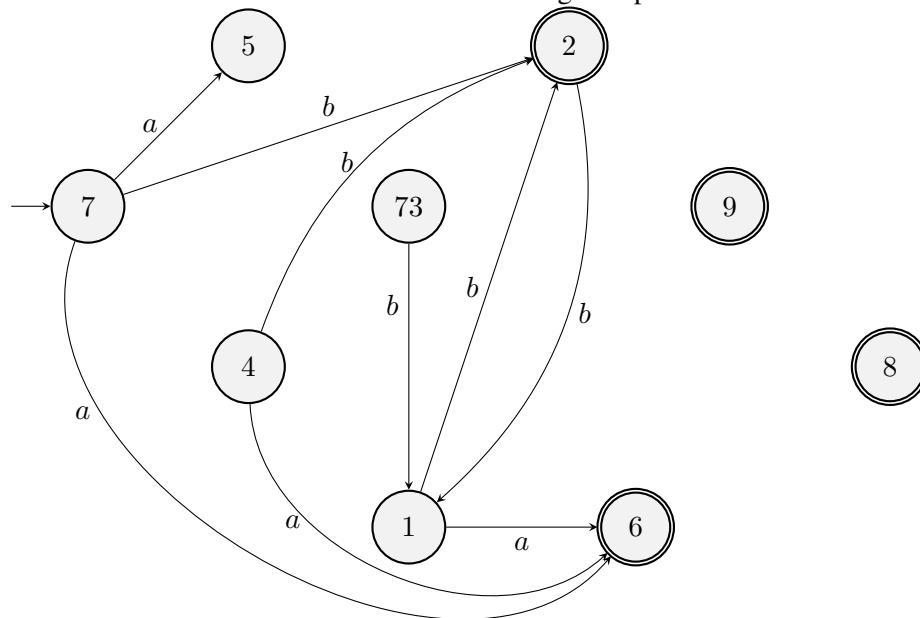
$2 \xrightarrow{\varepsilon} 73 \xrightarrow{b} 1$ is a *slow b-transition* from 2 to 1;

$1 \xrightarrow{b} 2$ is a *slow b-transition* from 1 to 2;

$2 \xrightarrow{\varepsilon} 9 \xrightarrow{\varepsilon} 8$ is *slowly accepting*.

You can see that this path consists of three *slow b-transitions* followed by *slow acceptance*. A *slow b-transition* consists of several (zero or more) ε -transitions, culminating in a *b-transition*. *Slow acceptance* consists of several ε -transitions, culminating in an accepting state.

Now let's see how to remove the ε -transitions to give a plain NFA.



The automaton looks similar to before: the states are the same and the initial state is the same. The difference is that the transitions you see here are the slow transitions, and the accepting states you see here are the slowly accepting states. (States 8, 9 and 4 can be removed because they're unreachable.)

6 Kleene's Theorem

See also the video "Kleene's Theorem" on Canvas (28 minutes).

Up to this point, we have seen examples of

- using a regexp to describe a language
- solving a matching problem on a DFA.

It turns out that these two things are closely connected.

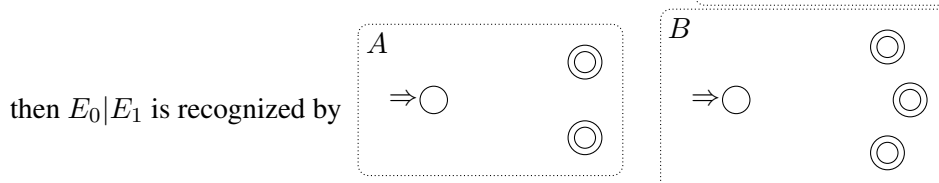
Theorem 1 (Kleene's Theorem) For a language $L \subseteq \Sigma^*$, the following are equivalent.

1. L is regular, i.e. it can be described by a regexp.
2. The matching problem for L can be solved by a DFA.

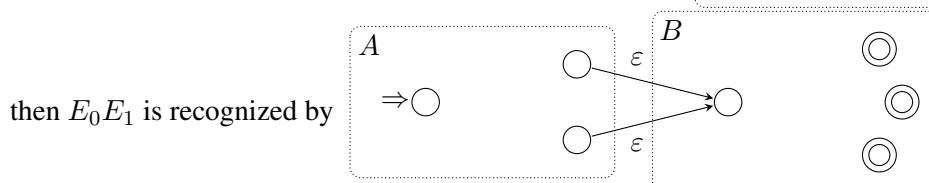
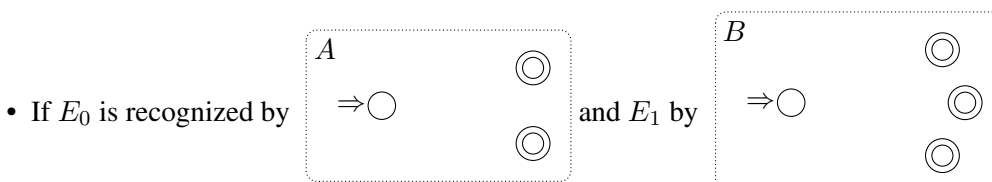
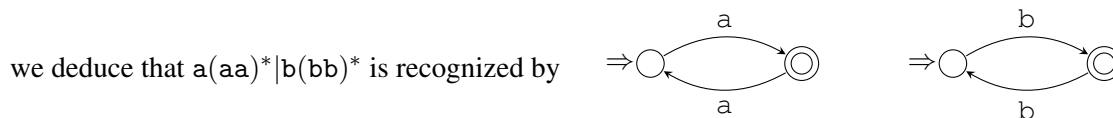
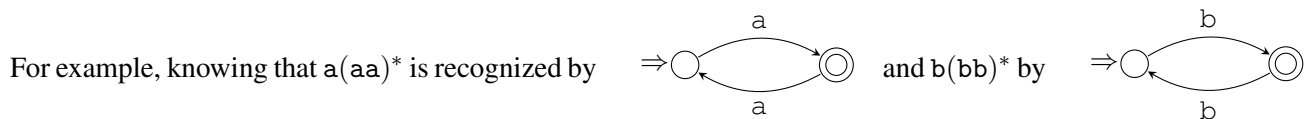
We shall only study the proof of $(1) \Rightarrow (2)$, i.e. we'll see how to turn a regexp into a DFA that recognizes the same language. As we saw, it suffices to construct an ε NFA, because then we remove the ε -transitions to obtain an NFA, and lastly we determinize to obtain a DFA.

So we want to convert a regexp into an ε NFA that recognizes the same language.

- The regexp a is recognized by $\Rightarrow \text{---} \xrightarrow{a} \text{---}$ and likewise if E is b or c .
- The regexp ε is recognized by $\Rightarrow \text{---}$



i.e. we take the sum of the two sets of states, and leave everything else unchanged.



For example, knowing that $a(aa)^*$ is recognized by $\Rightarrow \text{[Diagram: start state to loop state labeled 'a', then to final state]}$ and $b(bb)^*$ by $\Rightarrow \text{[Diagram: start state to loop state labeled 'b', then to final state]}$

we deduce that $a(aa)^*b(bb)^*$ is recognized by $\Rightarrow \text{[Diagram: start state to loop state labeled 'a', then epsilon transition to loop state labeled 'b', then to final state]}$

- If E is recognized by $\Rightarrow \text{[Diagram: start state to final state labeled 'A']}$ then E^* is recognized by $\Rightarrow \text{[Diagram: start state to loop state labeled 'A', then to final state]}$

For example, $(ab)^*a$ is recognized by $\Rightarrow \text{[Diagram: start state to loop state labeled 'a', then to final state]}$

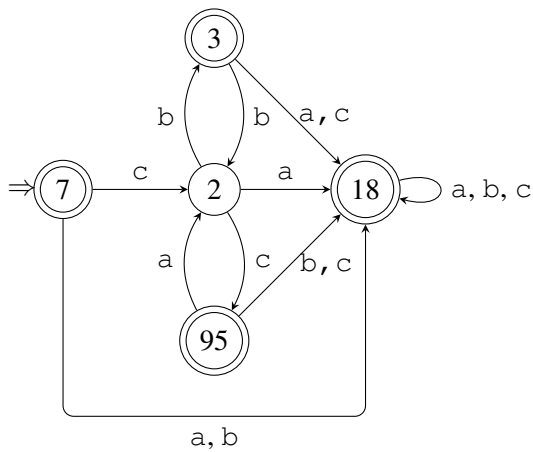
so $((ab)^*a)^*$ is recognized by $\Rightarrow \text{[Diagram: start state to loop state labeled 'a', then to final state]}$

- Finally, \emptyset is recognized by the empty automaton (the partial DFA with no states).

6.1 Using Kleene's Theorem

Although regexps and finite automata are equivalent, it's often easier to operate on automata. For example, we asked: is the complement of a regular language regular? That's not obvious for regexps, but using (total) DFAs, it is clear: just replace accepting states by non-accepting ones and vice versa. For example, for the alphabet $\{a, b, c\}$, we know

that $c(bb|ca)^*$ is recognized by $\Rightarrow \text{[Diagram: DFA with states 7, 2, 3, 95, 18. Transitions: 7 to 2 on 'c', 2 to 3 on 'b', 3 to 2 on 'b', 2 to 18 on 'a', 18 to 2 on 'a', 2 to 95 on 'c', 95 to 2 on 'c', 95 to 18 on 'b, c', 18 to 18 on 'a, b, c', 18 to 7 on 'a, b']}$ so its complement is recognized by



This implies that the intersection of regular languages L and M is regular, because $M \cap L = \overline{\overline{L} \cup \overline{M}}$. But we can also see this directly. Say L is recognized by automaton A , and M by B . If you're given a word w and you want to know whether it's in $L \cap M$, passing through w just once, you can simultaneously keep track of where you are in A and where you are in B , as you read through w . So your state at any time is a pair (x, y) , where x is a state in A , and y is a state in B . This immediately gives you a DFA.

For example, say you want a DFA for the password question in the last handout. You can make a DFA that determines whether a word has at least 3 characters, and another that determines whether it has a letter, and another that determines whether it has a digit. Then we want a DFA for the intersection of these languages. Kleene's theorem tells us that there must be some corresponding regexp.