

Tuesday - August 13, 2024

## Regular languages

### 1 Interchangeable subexpressions

Let's forget about FSAs for a moment, and just consider sets of strings "out of the blue." We'll connect things back to FSAs shortly. Here's a useful concept:

#### (1) L-remainders

- a. Given some stringset  $L$ , the  $L$ -remainders of a string  $u$  are all the strings  $v$  such that  $uv \in L$ . I'll write  $rem_L(u)$  for the set of  $L$ -remainders of  $u$ , so we can write this definition in symbols as:  $rem_L(u) = \{v \mid v \in \Sigma^*, uv \in L\}$ .

Roughly,  $rem_L(u)$  gives us a handle on "all the things we're still allowed to do, if we've done  $u$  so far."

- (2) If  $L_1 = \{\text{cat, cap, cape, cut, cup, dog}\}$ , then:

- a.  $rem_{L_1}(\text{ca}) = \{\text{t, p, e}\}$   
 b.  $rem_{L_1}(\text{c}) = \{\text{at, ap, ape, ut, up}\}$   
 c.  $rem_{L_1}(\text{cap}) = \{\text{e, t}\}$

ex:  $\{cv, cuv\}$   
 $rem_L(c) = \{v, vv\}$

empty string

- (3) If  $L_2 = \{\text{ad, add, baa, bad, cab, cad, dab, dad}\}$ , then:

- a.  $rem_{L_2}(\text{c}) = \{\text{ab, ad}\}$   
 b.  $rem_{L_2}(\text{d}) = \{\text{ab, ad}\}$   
 c.  $rem_{L_2}(\text{a}) = \{\text{d, dd}\}$   
 d.  $rem_{L_2}(\text{da}) = \{\text{b, d}\}$   
 e.  $rem_{L_2}(\text{ad}) = \{\text{b, d}\}$

empty string

When we notice that  $rem_{L_2}(\text{c}) = rem_{L_2}(\text{d})$ , this tells us something useful about how we can go about designing a grammar to generate the stringset  $L_2$ . Such a grammar doesn't need to care about the distinction between starting with 'c' and starting with 'd.' This is because, for any string  $v$  that you choose,  $cv$  and  $dv$  will either both be in  $L_2$  or both not be in  $L_2$ . An initial 'c' and an initial 'd' are thus interchangeable subexpressions.

- (4) **L-equivalence:** Given a stringset  $L$  and two strings  $u \in \Sigma^*$  and  $v \in \Sigma^*$ , we define a relation  $\equiv_L$  such that:  $u \equiv_L v$  iff  $rem_L(u) = rem_L(v)$ .

"(c) and (d) L equivalent since they have same L remainder"

Some slightly more linguistics-ish examples:

- (5) Suppose that  $\Sigma = \{C, V\}$ , and  $L$  contains all strings that contain at least one 'V'. Then:
- alphabet*
- L equivalent*
- $C \equiv_L CC$ , because both can only be followed by strings that fulfill the requirement for a 'V'.
  - $VC \equiv_L CV$ , because both can be followed by anything at all.
  - Two strings are  $L$ -equivalent iff they either both do or both don't contain a 'V'.
- (6) Suppose that  $\Sigma = \{C, V\}$ , and  $L$  contains all strings that have two adjacent 'C's or two adjacent 'V's (or both). Then:
- $C \equiv_L CVC \equiv_L CVCVC$ , because these all require remainders that have two adjacent 'C's, or two adjacent 'V's, or an initial 'C'. *create CC*
  - $V \equiv_L VCV \equiv_L VCVCV$ , because these all require remainders that have two adjacent 'C's, or two adjacent 'V's, or an initial 'V'. *create VV*
  - $CC \equiv_L VCVCVVCVC$  *both satisfy requirement already* *L remainder is any string in the alphabet for this ex*
- (7) Suppose that  $\Sigma = \{C, V\}$ , and  $L$  contains all strings that do not have two adjacent occurrences of 'V'. Then:
- $CCCC \equiv_L VC$ , because both can be followed by anything without adjacent 'V's.
  - $CCV \equiv_L V$ , because both can be followed by anything without adjacent 'V's that does not begin with 'V'. *don't want VV*
  - $CCV \not\equiv_L CCC$ , because only the latter can be followed by 'VC'.
  - So, two strings are  $L$ -equivalent if they end with the same symbol. *★*
- alphabet*
- (8) Suppose that  $\Sigma$  is the set of English words, and  $L$  is the set of all grammatical English-word sequences. Then (at least from the perspective of the syntax):
- John  $\equiv_L$  the brown furry rabbit *you can substitute these in to any english sentence*
  - John  $\equiv_L$  Mary thinks that John
  - John  $\not\equiv_L$  the fact that John *we can swap them*

## 2 The Myhill-Nerode Theorem

We can <sup>prefix or suffix</sup> connect this idea of equivalent subexpressions back to forward values in an FSA. Recall that  $fwd_M(u)(q)$  is a Boolean, True or False; so we can think of  $fwd_M(u)$  as a set of states, namely all those states reachable from an initial state of  $M$  by taking transitions that produce the string  $u$ . I'll sometimes call this a **forward set**, for lack of a better name.

Now here's the important connection:

- (9) For any FSA  $M = (Q, \Sigma, I, F, \Delta)$  and for any two strings  $u \in \Sigma^*$  and  $v \in \Sigma^*$ , if  $fwd_M(u) = fwd_M(v)$ , then  $u \equiv_{L(M)} v$ .  
<sup>both strings  $\perp$  equivalent if have same forward set</sup>

Given any particular stringset  $L$ , we can think of the relation  $\equiv_L$  as sorting out all possible strings into "buckets" (or **equivalence classes**): two strings belong in the same bucket iff they are equivalent prefixes.

- According to (9), for an FSA to generate  $L$  it must be arranged so that  $fwd$  only maps two strings to the same state-sets if those two strings are equivalent prefixes.
- The machine can ignore distinctions between bucket-mates, but *only* between bucket-mates.

And now we can put our finger on the capacities and limitations of finite-state automata:

- (10) **The Myhill-Nerode Theorem:** Given a particular stringset  $L$ , there is an FSA that generates  $L$  iff the relation  $\equiv_L$  sorts strings into only **finitely-many** buckets.

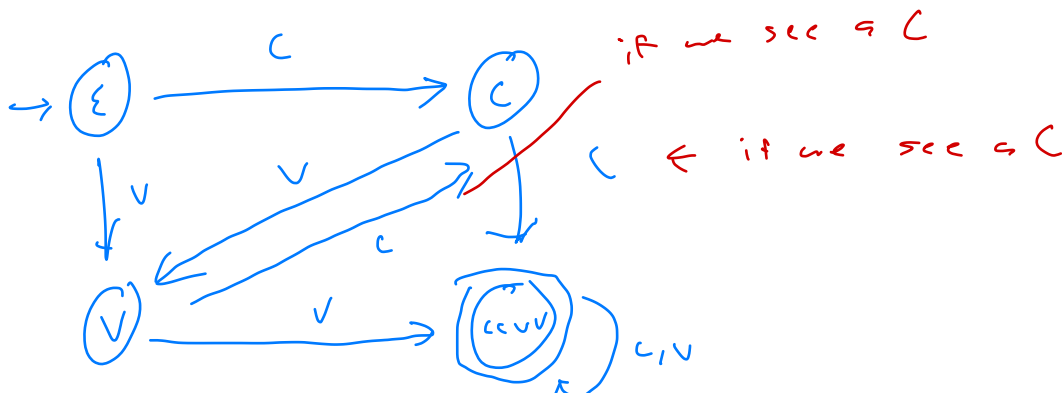
### 2.1 Why can any finitely-many distinctions be captured by an FSA?

If we have a particular stringset  $L$  whose equivalence relation  $\equiv_L$  makes only finitely-many distinctions, then there is a straightforward way to construct a minimal FSA whose states track exactly those distinctions.

Consider again the stringset from (6), consisting of all strings with either two adjacent 'C's or two adjacent 'V's (or both). <sup>6</sup> This stringset's equivalence relation sorts strings into four buckets:

- (11)
- a bucket containing strings that have either two adjacent 'C's or two adjacent 'V's;
  - a bucket containing strings that don't have two adjacent 'C's or 'V's, but end in 'C';
  - a bucket containing strings that don't have two adjacent 'C's or 'V's, but end in 'V';
  - a bucket containing only the empty string.

Having noticed this, we can mechanically construct an appropriate FSA— known as the **minimal FSA** for this stringset— which has one state corresponding to each bucket.



if two strings are equivalent,

The crucial idea here is that if  $u \equiv_L v$ , then  $ux \equiv_L vx$  for any  $x \in \Sigma$ . Adding a symbol at the end can't "break" an equivalence. Similarly, for any FSA  $M$ , if  $fwd_M(u) = fwd_M(v)$ , then  $fwd_M(ux) = fwd_M(vx)$ .

finite = finite if two strings forward set

## 2.2 Why do FSAs only make finitely-many distinctions?

On the other hand, if we have a particular FSA whose set of states is  $Q$ , then there are only finitely many distinct subsets of  $Q$  that  $fwd_M$  can map strings to, namely  $2^{|Q|}$  of them. This means that there are only finitely-many "candidate forward sets" for any string, so the FSA is necessarily making only those finitely-many distinctions between strings.

a partition of a

To illustrate this concept, consider the string set  $L = \{a^n b^n \mid n > 0\}$ , which no FSA is able to generate. Notice that  $a \not\equiv_L aa$ , and  $aa \not\equiv_L aaa$ , and so on.

not regular

- In fact, any string of 'a's is non-equivalent to any other different-length string of 'a's, so this stringset sorts strings into infinitely-many buckets, one bucket for each length.
- But there is no way for an FSA  $M$  to be set up such that  $fwd_M(a^j) \not\equiv_L fwd_M(a^k)$  whenever  $j \neq k$ ; any FSA will incorrectly collapse the distinction between two such strings of 'a's.

Regular language = any language captured by FSA

### 3 FSAs with epsilon transitions

Epsilon transitions provide another useful way to extend standard FSAs:

(12) A finite-state automaton with epsilon transitions ( $\epsilon$ -FSA) is a five-tuple  $(Q, \Sigma, I, F, \Delta)$  where:

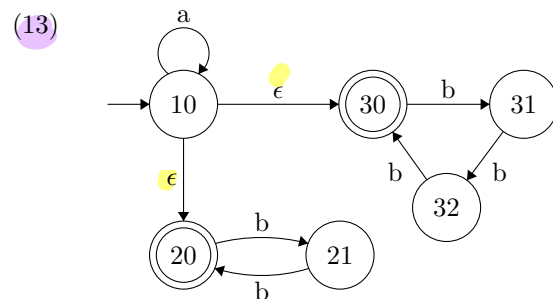
- $Q$  is a finite set of states;
- $\Sigma$ , the alphabet, is a finite set of symbols;
- $I \subseteq Q$  is the set of initial states;
- $F \subseteq Q$  is the set of ending states;
- $\Delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times Q$  is the set of transitions.

states  
alphabet  
initial  
final

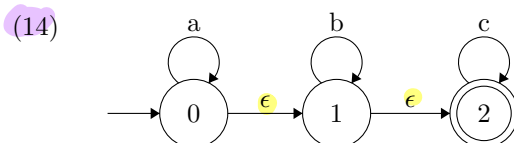
transition over empty string

The difference from what we've seen before is that transitions are labeled not with an element of  $\Sigma$ , but rather with *either* an element of  $\Sigma$  *or* the empty string. (Note that  $\epsilon$  is a member of the set  $\Sigma^*$ , but it is not a member of  $\Sigma$ .)

Here are two  $\epsilon$ -FSAs for illustration:



epsilon closure of 10  
 $\{10, 20, 30\}$



The key idea in dealing with  $\epsilon$ -transitions is the **epsilon closure** of a state:

ex: epsilon closure of 0  
 $\{0, 1, 2\}$

(15) If  $\Delta$  is the transition function,  $cl_{\Delta}(q)$  is the set of all states reachable from  $q$  by a sequence of zero-or-more  $\epsilon$ -transitions according to  $\Delta$ .

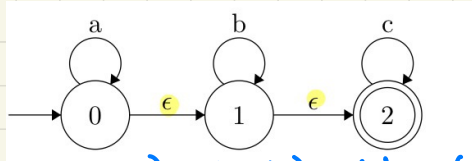
It turns out that an FSA with epsilon transitions can always be converted into another FSA that does not contain any epsilon transitions, but generates exactly the same stringset.

Here's the recipe for converting an  $\epsilon$ -FSA into an FSA that does not contain epsilon-transitions:

(16) Given an  $\epsilon$ -FSA  $M = (Q, \Sigma, I, F, \Delta)$  (which may contain epsilon transitions), we can construct an FSA  $M' = (Q, \Sigma, I, F', \Delta')$  (which does not contain epsilon transitions) that will generate the same stringset as  $M$  as follows:

- The new set of end states,  $F'$ , contains all states  $q$  such that  $cl_{\Delta}(q) \cap F \neq \emptyset$
- The new transition set  $\Delta'$  contains a transition  $(q_1, x, q_3)$  iff there is some  $q_2 \in cl_{\Delta}(q_1)$  such that  $(q_2, x, q_3) \in \Delta$

ex using 14



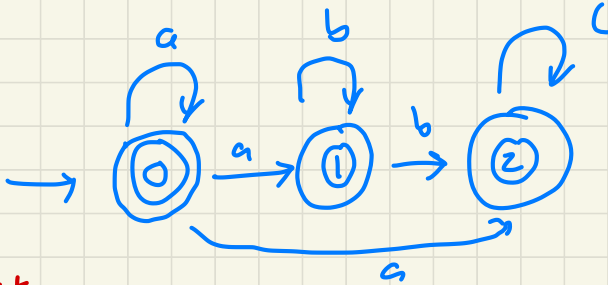
step 1:

$$\begin{aligned}
 & \mathcal{L}(\Delta(0)) = \{0, 1, 2\} & \mathcal{L}(\Delta(1)) = \{1, 2\} & \mathcal{L}(\Delta(2)) = \{2\} \\
 & \cap & \cap & \cap \\
 & = \{2\} & = \{2\} & = \{2\}
 \end{aligned}$$

"look  
for  
final  
states

$$F = \{2\}$$

↗ we can  
create  
 $a^*b^*c^*$



we can get  
from 0 to 1 and 2  
by just using  
a

we can get from  
1 to 2 just using  
b

## 4 Regular expressions and FSAs

### 4.1 Defining regular expressions and their denotations

First we'll define what regular expressions are. That's all we're saying in (17). It's analogous to defining what counts as a propositional formula.

(17) Given an alphabet  $\Sigma$ , we define  $\text{RE}(\Sigma)$ , the set of regular expressions over  $\Sigma$ , as follows:

- $\{ \}$  / regular expression
- $\mathbf{0} \in \text{RE}(\Sigma)$
  - $\mathbf{1} \in \text{RE}(\Sigma)$
  - if  $x \in \Sigma$ , then  $\underline{x} \in \text{RE}(\Sigma)$
  - if  $\underline{r_1} \in \text{RE}(\Sigma)$  and  $\underline{r_2} \in \text{RE}(\Sigma)$ , then  $\underline{(r_1 \mid r_2)} \in \text{RE}(\Sigma)$  or
  - if  $\underline{r_1} \in \text{RE}(\Sigma)$  and  $\underline{r_2} \in \text{RE}(\Sigma)$ , then  $\underline{(r_1 \cdot r_2)} \in \text{RE}(\Sigma)$  concatenated
  - if  $\underline{r} \in \text{RE}(\Sigma)$ , then  $\underline{r^*} \in \text{RE}(\Sigma)$

So if we have the alphabet  $\Sigma = \{a, b, c\}$ , then here are some elements of  $\text{RE}(\Sigma)$ :

- $\underline{(a \mid b)}$
  - $\underline{((a \mid b) \cdot c)}$
  - $\underline{((a \mid b) \cdot c)^*}$
- $\leftarrow$  in haskell code example

Now, any regular expression  $r \in \text{RE}(\Sigma)$  denotes a stringset. We'll write  $\llbracket r \rrbracket$  for the stringset denoted by  $r$ .

(18) Given a regular expression  $r \in \text{RE}(\Sigma)$ ,  $\llbracket r \rrbracket \subseteq \Sigma^*$  such that:

- $\llbracket \mathbf{0} \rrbracket = \emptyset = \{ \}$
  - $\llbracket \mathbf{1} \rrbracket = \{ \epsilon \}$
  - $\llbracket \underline{x} \rrbracket = \{ x \}$
  - $\llbracket \underline{(r_1 \mid r_2)} \rrbracket = \llbracket r_1 \rrbracket \cup \llbracket r_2 \rrbracket$  union
  - $\llbracket \underline{(r_1 \cdot r_2)} \rrbracket = \{ u ++ v \mid u \in \llbracket r_1 \rrbracket, v \in \llbracket r_2 \rrbracket \}$  concatenate
  - $\llbracket \underline{r^*} \rrbracket$  is the smallest set such that:
    - $\epsilon \in \llbracket r^* \rrbracket$
    - if  $u \in \llbracket r \rrbracket$  and  $v \in \llbracket r^* \rrbracket$ , then  $u ++ v \in \llbracket r^* \rrbracket$
- $\} \text{ zero or more repeats of } r$

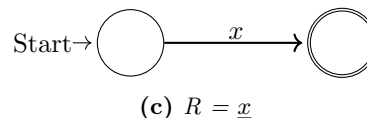
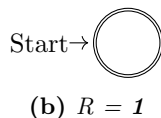
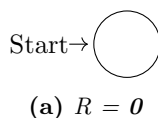
The tricky part here is the  $r^*$  case. It says roughly that  $\llbracket r^* \rrbracket$  is the set comprising all strings that we can get by concatenating zero or more strings from the set  $\llbracket r \rrbracket$ .

- Concatenating *zero* such strings produces  $\epsilon$ , so  $\epsilon \in \llbracket r^* \rrbracket$ .
- Concatenating *n* such strings, where *n* is *non-zero*, really means concatenating some string  $u$ , which is in  $\llbracket r \rrbracket$ , with some string  $v$ , which is the result of concatenating some  $n - 1$  such strings.

## 4.2 Relating regular expressions to FSAs

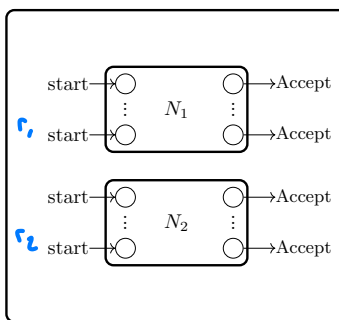
It turns out that given any regular expression  $r$ , we can construct an  $\epsilon$ -FSA  $M$  such that  $\mathcal{L}(M) = \llbracket r \rrbracket$ . That is, we can construct an  $\epsilon$ -FSA that generates exactly the stringset denoted by  $r$ . To do this, we have to proceed recursively on the structure of the regular expression, because there are unboundedly many regular expressions to deal with. The diagrams below give the important idea for how this works.

- ...  $\llbracket 0 \rrbracket$ ,  $\llbracket 1 \rrbracket$  and  $\llbracket x \rrbracket$ , for any  $x \in \Sigma$ , as follows



*also add state*

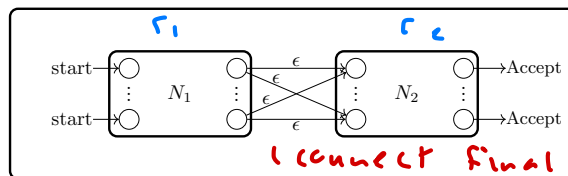
- ...  $\llbracket (r_1 \mid r_2) \rrbracket$ , given an  $\epsilon$ -FSA  $N_1$  that generates  $\llbracket r_1 \rrbracket$  and an  $\epsilon$ -FSA  $N_2$  that generates  $\llbracket r_2 \rrbracket$ , as follows:



*yes*  
 $\leftarrow$  new fsa has everything from both

Figure 2: The construction used in the union of two  $\epsilon$ -FSAs

- ...  $\llbracket (r_1 \cdot r_2) \rrbracket$ , given an  $\epsilon$ -FSA  $N_1$  that generates  $\llbracket r_1 \rrbracket$  and an  $\epsilon$ -FSA  $N_2$  that generates  $\llbracket r_2 \rrbracket$ , as follows.

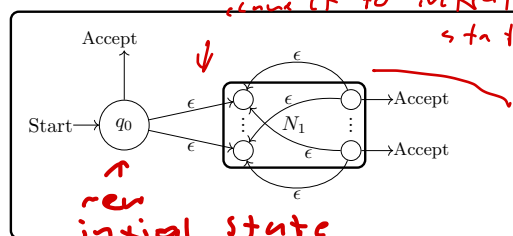


*concatenated*

*connect final states of  $r_1$  with starting states in  $r_2$  using epsilon transitions*

Figure 3: The construction used in the concatenation of two  $\epsilon$ -FSAs

- ...  $\llbracket r^* \rrbracket$ , given an  $\epsilon$ -FSA  $N_1$  that generates  $\llbracket r \rrbracket$ .



*connect final to initial states using epsilon*

Figure 4: The construction used in the star operation



## 5 Summing up: Regular languages

A stringset  $S \subseteq \Sigma^*$  is a **regular stringset** (or **regular language**) iff there is some regular expression  $r \in \text{RE}(\Sigma)$  such that  $\llbracket r \rrbracket = S$ .

*language of FSA*

We've just seen that, given any regular expression  $r$ , we can construct an  $\epsilon$ -FSA  $M$  such that  $\mathcal{L}(M) = \llbracket r \rrbracket$ . This tells us that any stringset that can be described by a regular expression — any regular language — can also be described by an  $\epsilon$ -FSA.

It turns out that the inverse is also true: given any  $\epsilon$ -FSA  $M$ , there is a method for constructing a regular expression  $r$  such that  $\llbracket r \rrbracket = \mathcal{L}(M)$ . So this means that any stringset that can be described by an  $\epsilon$ -FSA is a regular language.<sup>1</sup>

So in summary, for any stringset  $S$ , either all of the following are true or all are false:

- (19)
- a.  $S$  is a regular language.
  - b. There is some regular expression  $r$  such that  $\llbracket r \rrbracket = S$ .
  - c. There is some  $\epsilon$ -FSA  $M$  such that  $\mathcal{L}(M) = S$ .
  - d. There is some FSA  $M$  such that  $\mathcal{L}(M) = S$ .
  - e. The number of distinct classes of equivalent prefixes, i.e. the number of buckets into which strings are sorted by the relation  $\equiv_L$ , is finite.
  - f. The number of distinct classes of equivalent suffixes is finite.<sup>2</sup>

*one of things is true*

<sup>1</sup>This is trickier to prove, but see e.g. pp. 33-34 Hopcroft and Ullman (1979) or pp. 67-74 of Sipser (1997).

<sup>2</sup>This is just the obvious parallel to the idea of equivalent prefixes: two strings  $u$  and  $v$  are equivalent suffixes iff, for all strings  $w$ ,  $wu$  and  $wv$  are either both in  $S$  or both not in  $S$ .