

### 3. Introducing finite-state automata

First some standard stage-setting definitions:

- (1) For any set  $\Sigma$ , we define  $\Sigma^*$  as the smallest set such that:
- $\epsilon \in \Sigma^*$ , and
  - if  $x \in \Sigma$  and  $u \in \Sigma^*$  then  $(x:u) \in \Sigma^*$ .

We often call  $\Sigma$  an *alphabet*, call the members of  $\Sigma$  *symbols*, and call the members of  $\Sigma^*$  *strings*.

- (2) For any two strings  $u \in \Sigma^*$  and  $v \in \Sigma^*$ , we define  $u \mathbin{++} v$  as follows:
- $\epsilon \mathbin{++} v = v$
  - $(x:w) \mathbin{++} v = x:(w \mathbin{++} v)$

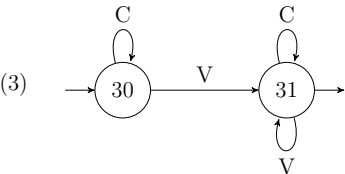
Although these definitions provide the “official” notation, I’ll sometimes be slightly lazy and abbreviate ‘ $x:\epsilon$ ’ as ‘ $x$ ’, and abbreviate both ‘ $s:t$ ’ and ‘ $s \mathbin{++} t$ ’ as just ‘ $st$ ’ in cases where it should be clear what’s intended.

I’ll generally use  $x$ ,  $y$  and  $z$  for individual symbols of an alphabet  $\Sigma$ , and use  $u$ ,  $v$  and  $w$  for strings in  $\Sigma^*$ . This should help to clarify whether a ‘:’ or a ‘ $\mathbin{++}$ ’ has been left out.

#### 1 Finite-state automata, informally

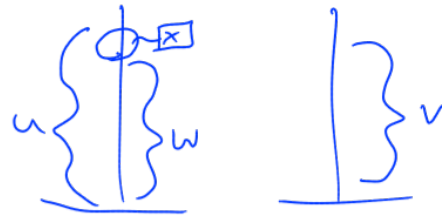
Below, in (3) and (4), are graphical representations of two finite-state automata (FSAs). The circles represent *states*. The *initial* states are indicated by an “arrow from nowhere”; the *final* or *accepting* states are indicated by an “arrow to nowhere”.

The FSA in (3) generates the subset of  $\{C, V\}^*$  consisting of all and only strings that have at least one occurrence of ‘V’.



The FSA in (4) generates the subset of  $\{C, V\}^*$  consisting of all and only strings that contain either two adjacent ‘C’s or two adjacent ‘V’s (or both).

$$\Sigma = \{a, b, c\}$$
$$\Sigma^* = \{aab, bccba, \dots\}$$



concatenate ::  $[a] \rightarrow [a] \rightarrow [a]$

concatenate =  $\lambda u \rightarrow \lambda v \rightarrow$

$[ ] \rightarrow v$

$x:w \rightarrow x:(\text{concatenate } w \ v)$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

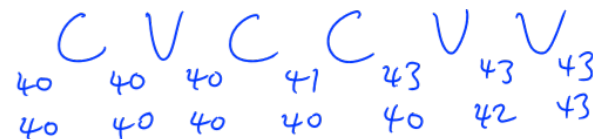
$$\Sigma = \{\text{cat, dog, the, big, } \dots\}$$

$$\Sigma = \{p, k, t, s, z, j, \dots\}$$

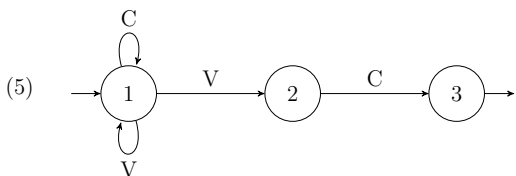
$$\Sigma = \{C, V\}$$

$$\Sigma^* = \{k\check{a}et, bl\check{i}k, bu\check{u}k, \dots\}$$

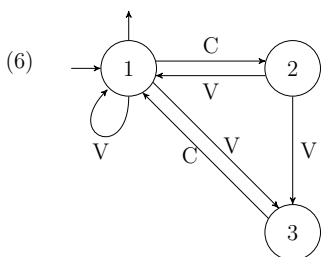
$$\Sigma^* = \{V, CV, VC, CVC, CVCC, CCVC, VVC, VVVV, \dots\}$$



A.	CCV	✓
B.	VVVVVV	✓
C.	VCV	X
D.	VCCV	✓
E.	VCVC	X



# Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

- $\{c, v\}$

- $Q$  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; and
- $\Delta \subseteq Q \times \Sigma \times Q$  is the set of transitions.

(8) ( {40, 41, 42, 43}, {C, V}, {40}, {43},  
{ (40, C, 40), (40, C, 41), (40, V, 40), (40, V, 42), (41, C, 43), (42, V, 43), (43, C, 43), (43, V, 43) } )

Now let's try to say more precisely what it means for an automaton  $M = (Q, \Sigma, I, F, \Delta)$  to generate/accept a string.

- (9) For  $M$  to generate a string of three symbols, say  $x_1x_2x_3$ , there must be four states  $q_0, q_1, q_2$ , and  $q_3$  such that

- $q_0 \in I$ , and
- $(q_0, x_1, q_1) \in \Delta$ , and
- $(q_1, x_2, q_2) \in \Delta$ , and
- $(q_2, x_3, q_3) \in \Delta$ , and
- $q_3 \in F$ .

$$q_0 \xrightarrow{x_1} q_1 \xrightarrow{x_2} q_2 \xrightarrow{x_3} q_3$$

- (10) More generally,  $M$  generates a string of  $n$  symbols, say  $x_1 x_2 \dots x_n$ , iff: there are  $n + 1$  states  $q_0, q_1, q_2, \dots, q_n$  such that
- $q_0 \in I$ , and
  - for every  $i \in \{1, 2, \dots, n\}$ ,  $(q_{i-1}, x_i, q_i) \in \Delta$ , and
  - $q_n \in F$ .

$$q_0 \xrightarrow{x_1} q_1 \xrightarrow{x_2} q_2 \dots q_{i-1} \xrightarrow{x_i} q_i \dots q_{n-1} \xrightarrow{x_n} q_n$$

$|Q|^{n+1}$  possible state-sequences

To take a concrete example:

- (11) The automaton in (4)/(8) generates the string 'VCCVC' because we can choose  $q_0, q_1, q_2, q_3, q_4$  and  $q_5$  to be the states 40, 40, 41, 43, 43 and 43 (respectively), and then it's true that:
- $40 \in I$ , and
  - $(40, V, 40) \in \Delta$ , and
  - $(40, C, 41) \in \Delta$ , and
  - $(41, C, 43) \in \Delta$ , and
  - $(43, V, 43) \in \Delta$ , and
  - $(43, C, 43) \in \Delta$ , and
  - $43 \in F$ .

## Assignment Project Exam Help

**Side remark:** Note that abstractly, (10) is not all that different from:

- (12) A tree-based grammar will generate a string  $x_1 x_2 \dots x_n$  iff: there is some collection of nonterminal symbols that we can choose such that
- those nonterminal symbols and the symbols  $x_1, x_2$ , etc. can all be clicked together into a tree structure in ways that the grammar allows, and
  - the nonterminal "at the top" is the start symbol.

(Much more on this in a few weeks!)

<https://powcoder.com>

Add WeChat powcoder

$$\sum_{x \in \{1, 2, 3\}} (x^2) = 1^2 + 2^2 + 3^2$$

$$\bigvee_{x \in \{1, 2, 3\}} (x < 3) = (1 < 3) \vee (2 < 3) \vee (3 < 3) = \text{true}$$

We'll write  $\mathcal{L}(M)$  for the set of strings generated by an FSA  $M$ . So stated roughly, the important idea is:

- (13)  $w \in \mathcal{L}(M)$

$$\iff \bigvee_{\text{all possible paths } p} [\text{string } w \text{ can be generated by path } p]$$

$$\iff \bigvee_{\text{all possible paths } p} \left[ \bigwedge_{\text{all steps } s \text{ in } p} [\text{step } s \text{ is allowed and generates the appropriate part of } w] \right]$$

$\exists p [w \text{ can be generated by path } p]$

It's handy to write  $I(q_0)$  in place of  $q_0 \in I$ , and likewise for  $F$  and  $\Delta$ . Then one way to make (13) precise is:

- (14)  $x_1 x_2 \dots x_n \in \mathcal{L}(M)$

$$\iff \bigvee_{q_0 \in Q} \bigvee_{q_1 \in Q} \dots \bigvee_{q_{n-1} \in Q} \bigvee_{q_n \in Q} [I(q_0) \wedge \Delta(q_0, x_1, q_1) \wedge \dots \wedge \Delta(q_{n-1}, x_n, q_n) \wedge F(q_n)]$$

But it's practical and enlightening to break this down in a couple of different ways.

2.1 Forward values

For any FSA  $M$  there's a two-place predicate  $\text{fwd}_M$ , relating states to strings in an important way:

(15)  $\text{fwd}_M(w)(q)$  is true iff there's a path through  $M$  from some initial state to the state  $q$ , emitting the string  $w$

Given a way to work out  $\text{fwd}_M(w)(q)$  for any string and any state, we can easily use this to check for membership in  $\mathcal{L}(M)$ :

(16) 
$$w \in \mathcal{L}(M) \iff \bigvee_{q_n \in Q} [\text{fwd}_M(w)(q_n) \wedge F(q_n)]$$

We can represent the predicate  $\text{fwd}_M$  in a table. Each column shows  $\text{fwd}_M$  values for the *entire prefix* consisting of the header symbols to its *left*. The first column shows values for the empty string.

(17) Here's the table of  $\text{fwd}_M$  values for prefixes of the string 'CVCCVVC' for the FSA in (5).

State	C	V	C	C	V	V	C
1	1	1	1	1	1	1	1
2	0	0	1	0	0	1	0
3	0	0	0	1	0	0	1

Notice that filling in the values in the leftmost column is easy: this column just says which states are initial states. And with a little bit of thought you should be able to convince yourself that, in order to fill in a column of this table, you only need to know:

- the values in the column immediately to its left, and
- the symbol immediately to its left.

More generally, this means that:

(18) The  $\text{fwd}_M$  values for a non-empty string  $x_1 \dots x_n$  depend only on

- the  $\text{fwd}_M$  values for the string  $x_1 \dots x_{n-1}$ , and
- the symbol  $x_n$ .

This means that we can give a recursive definition of  $\text{fwd}_M$ :

(19) 
$$\begin{aligned} \text{fwd}_M(\epsilon)(q) &= I(q) \\ \text{fwd}_M(x_1 \dots x_n)(q) &= \bigvee_{q_{n-1} \in Q} [\text{fwd}_M(x_1 \dots x_{n-1})(q_{n-1}) \wedge \Delta(q_{n-1}, x_n, q)] \end{aligned}$$

This suggests a natural and efficient algorithm for calculating these values: write out the table, start by filling in the leftmost column, and then fill in other columns from left to right. This is where the name "forward" comes from.

2.2 Backward values

We can do all the same things, flipped around in the other direction.

For any FSA  $M$  there's a two-place predicate  $\text{bwd}_M$ , relating states to strings in an important way:

(20)  $\text{bwd}_M(w)(q)$  is true iff there's a path through  $M$  from the state  $q$  to some ending state, emitting the string  $w$

$$\text{fwd}(w)(q)$$

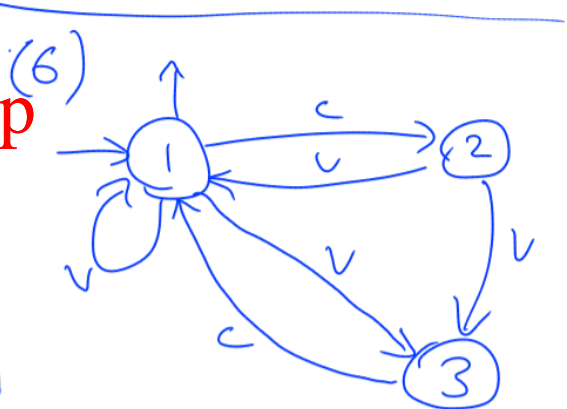


$$w = x_1 x_2 \dots x_{n-1} x_n$$
  
$$q_0 \quad q_1 \quad q_2 \quad \dots \quad q_{n-1} \quad q_n$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



	C	V	C	C	V	C
1					1	1
2					0	1
3					1	0

$$\text{fwd}(CVCCVC)(1) ?$$
  
$$\exists q' [\text{fwd}(CVCCV)(q') \wedge \Delta(q', C, 1)] ?$$

Given a way to work out  $\text{bwd}_M(w)(q)$  for any string and any state, we can easily use this to check for membership in  $\mathcal{L}(M)$ :

(21) 
$$w \in \mathcal{L}(M) \iff \bigvee_{q_0 \in Q} [I(q_0) \wedge \text{bwd}_M(w)(q_0)]$$

We can represent the predicate  $\text{bwd}_M$  in a table. Each column shows  $\text{bwd}_M$  values for the *entire suffix* consisting of the header symbols to its *right*. The last column shows values for the empty string.

(22) Here's the table of  $\text{bwd}_M$  values for suffixes of the string 'CVCCVVC' for the FSA in (5).

State	C	V	C	C	V	V	C
1	1	1	1	1	1	1	0
2	0	0	0	0	0	0	1
3	0	0	0	0	0	0	0

In this case, filling in the last column is easy, and each other column can be filled in simply by looking at the values immediately to its right.

(23) The  $\text{bwd}_M$  values for a non-empty string  $x_1 \dots x_n$  depend only on

- the  $\text{bwd}_M$  values for the string  $x_2 \dots x_n$ , and
- the symbol  $x_1$ .

So  $\text{bwd}_M$  can also be defined recursively.

(24) 
$$\begin{aligned} \text{bwd}_M(\epsilon)(q) &= F(q) \\ \text{bwd}_M(x_1 \dots x_n)(q) &= \bigvee_{q_1 \in Q} [\Delta(q, x_1, q_1) \wedge \text{bwd}_M(x_2 \dots x_n)(q_1)] \end{aligned}$$

2.3 Forward values and backward values together

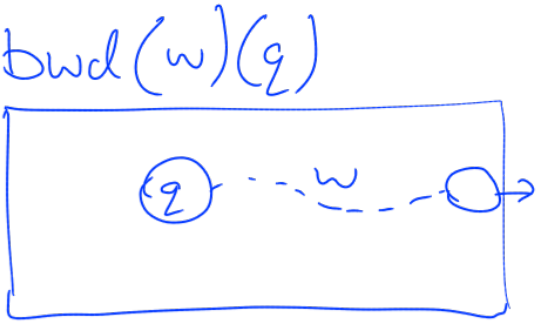
Now we can say something beautiful:

(25) 
$$uv \in \mathcal{L}(M) \iff \bigvee_{q \in Q} [\text{fwd}_M(u)(q) \wedge \text{bwd}_M(v)(q)]$$

And in fact (16) and (21) are just special cases of (25), with  $u$  or  $v$  chosen to be the empty string:

(26) 
$$w \in \mathcal{L}(M) \iff \bigvee_{q \in Q} [\text{fwd}_M(w)(q) \wedge \text{bwd}_M(\epsilon)(q)] \iff \bigvee_{q \in Q} [\text{fwd}_M(w)(q) \wedge F(q)]$$

(27) 
$$w \in \mathcal{L}(M) \iff \bigvee_{q \in Q} [\text{fwd}_M(\epsilon)(q) \wedge \text{bwd}_M(w)(q)] \iff \bigvee_{q \in Q} [I(q) \wedge \text{bwd}_M(w)(q)]$$



$$\text{bwd}(vvc)(1) = 1$$

$$\text{bwd}(vc)(3) = 0$$

$$\text{bwd}(\epsilon)(3) = 1$$

<https://powcoder.com>

Add WeChat powcoder

$$\text{bwd}(cvccvvc)(1) ?$$
  
$$\exists q_1. [\Delta(1, c, q_1) \wedge \text{bwd}(vccvvc)(q_1)] ?$$

