

Ling/CSE 472: Introduction to Computational Linguistics

3/28/12

Regular expressions & Finite State Automata

Overview

- Twitter: tweet classes
- Formal definition of regular languages
- Regular expression examples
- Extension to regular expressions
- Substitution in Python
- Finite state automata
- Recognition with FSA -- pseudocode
- DFSA v. NFSA

Announcement

- Assignment 1 is posted
- Elizalike.py (version 1) due 4/3 at 10am
- Complete assignment due 4/6 at 5pm

@everyone: Tweet the class! #cl472

- Goals
 - Create a sense of community within the class
 - Provide on-going feedback to peers and instructors about what is interesting/difficult/etc in the material
 - Understand the language of tweets
 - Create a (labeled) corpus for studying from the perspective of NLP
 - Digital literacy...
- What subcategory hashtags should we use?

Possible classes/thing we might do with tweets

- Emotion trying to convey: #wtf, #sarcasm, ...
- Parts of the class: #hw, #lecture, #reading, #general
- CS v. Ling: #cs, #ling
- Subtopic/subtask: #pos, #lemmatization, #...
- Dialogue acts: #question, #response, #fyi, #applications
- Formality level: #formal, #informal
- Idiomatic expressions

Hashtags: emotions + speech acts

- #curious, #help, #fyi, #wtf, #hth, #argh,

Formal languages

- From the point of view of formal languages theory, a language is a set of strings defined over some alphabet.
- Chomsky hierarchy: a description of classes of formal languages:
 - regular < context-free < context-sensitive < all languages
- Languages from a single level in the hierarchy can be described in terms of the same formal devices.

Three views on the same object

Regular languages can be described by regular expressions and by finite-state automata.

- Regular language: a set of strings
- Regular expression: an expression from a certain formal language which describes a regular language
- Finite-state automaton: a simple computing machine which accepts or generates a regular language

Formal definition of regular languages: symbols

- ϵ is the empty string
- \emptyset is the empty set
- Σ is an alphabet (set of symbols)
- Examples of alphabets:
 - $\Sigma = \{a,b,! \}$
 - $\Sigma = \{c,a,t,d,o,g,f,o,x \}$
 - $\Sigma = \{\text{cat, dog, fox} \}$
 - $\Sigma = \{1, 3, 5, 7 \}$

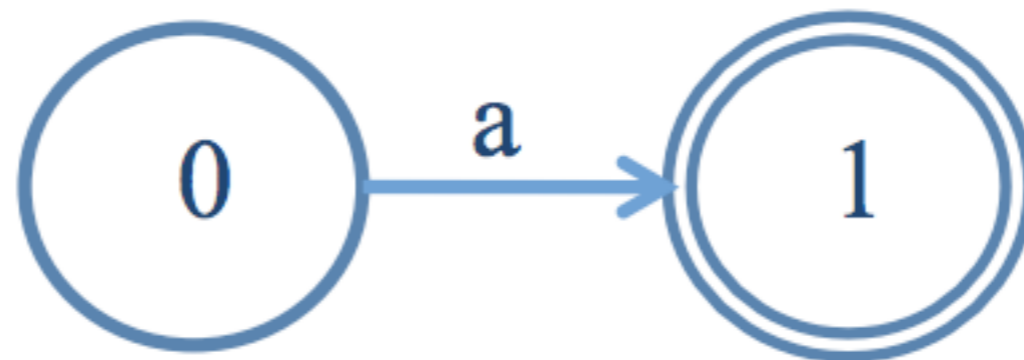
Formal definition of regular languages

The class of regular languages over Σ is formally defined as:

- \emptyset is a regular language
- $\forall a \in \Sigma \cup \epsilon, \{a\}$ is a regular language
- If $L1$ and $L2$ are regular languages, then so are:
 - $L1 \cdot L2 = \{xy \mid x \in L1, y \in L2\}$ (concatenation)
 - $L1 \cup L2$ (union or disjunction)
 - $L1^*$ (Kleene closure)

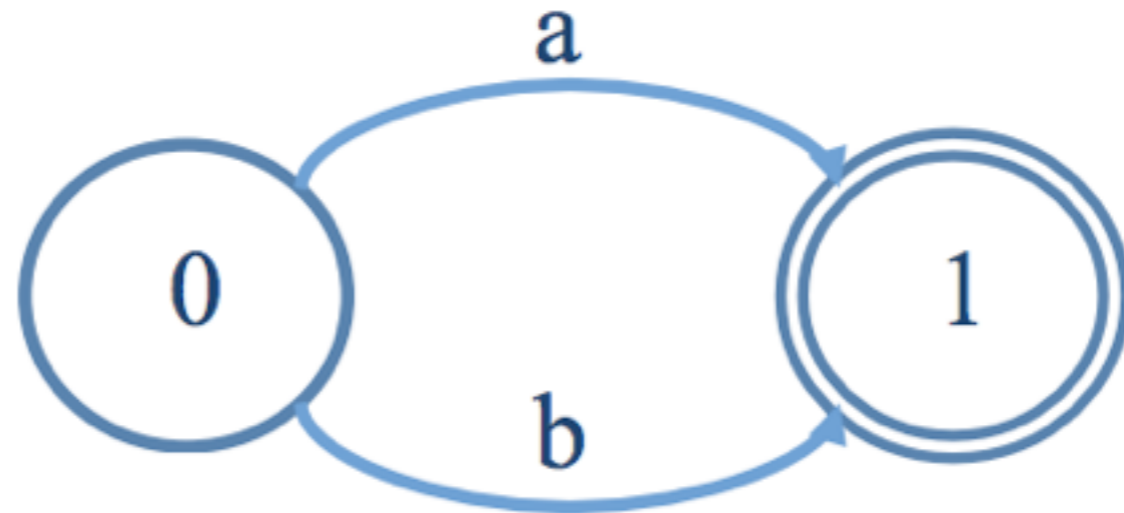
Examples: Single Character

- Alphabet: $\Sigma = \{a\}$
- Language (set of strings): $L = \{a\}$
- Regular expression: $/a/$
- Finite state machine:



Examples: Disjunction/union

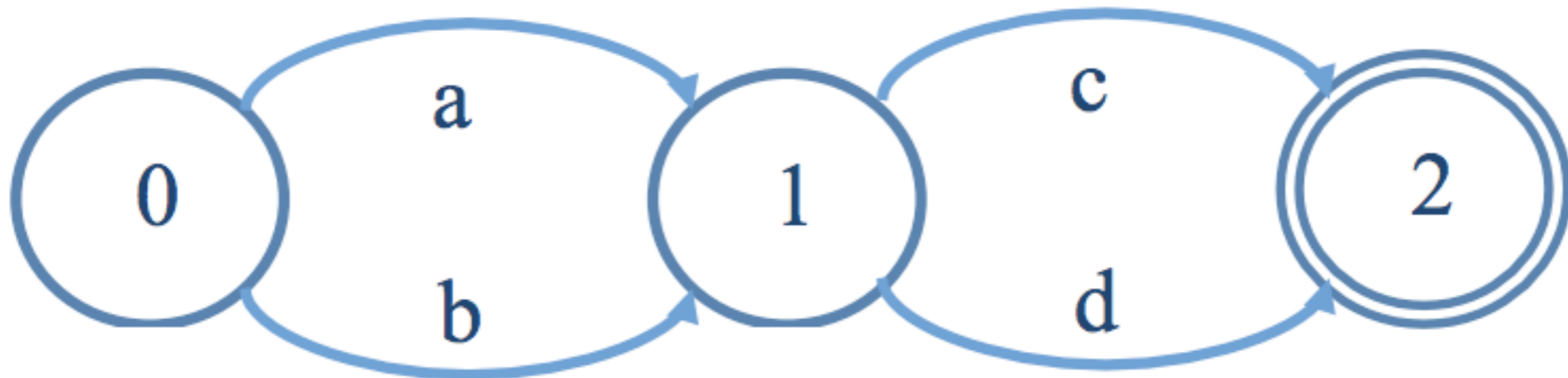
- Alphabet: $\Sigma = \{a, b\}$
- Language (set of strings): $\{a, b\}$
- Regular expression: $/(a|b)/$
- Finite state machine:



- This language/regexp/FSA also reflects the union of two languages: $\{a\} \cup \{b\}$

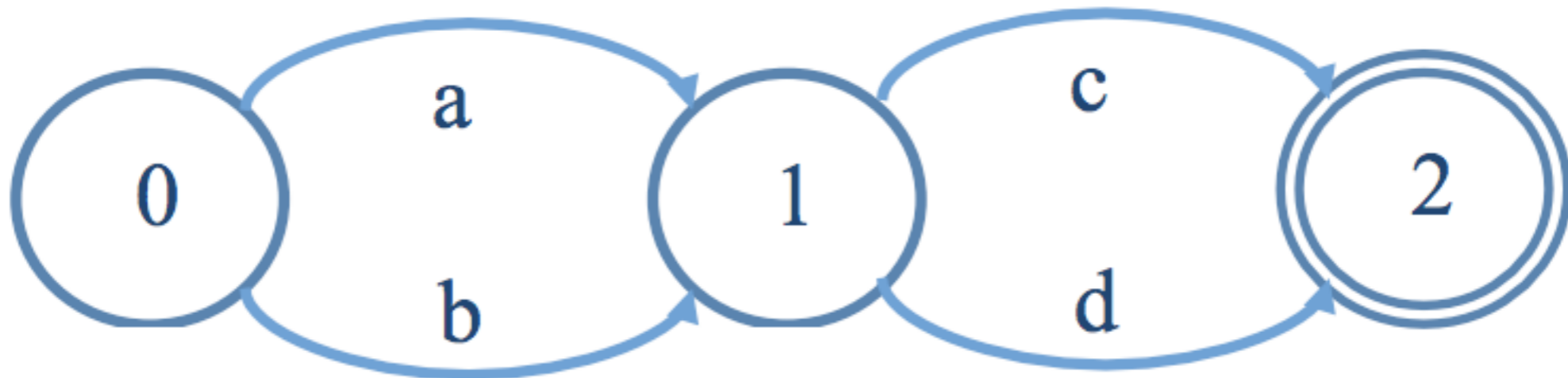
Examples: Concatenation

- $\Sigma = \{a, b, c, d\}$
- $L = \{ac, ad, bc, bd\}$
- Regular expression: $/(a|b)(c|d)/$
- Finite state machine:



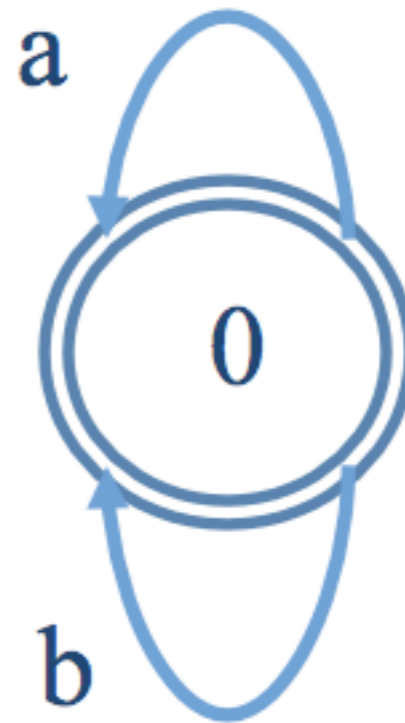
Examples: Concatenation

- Note: This language/regexp/FSA also reflects the concatenation of two languages:
 - $L1 = \{a, b\} \cdot L2 = \{c, d\}$



Examples: Kleene closure

- $\Sigma = \{a, b\}$ (again)
- $L = \{\epsilon, a, b, aa, bb, ab, ba, bb, \dots\}$
- Regular expression: $/(a|b)^*/$
- Finite state machine:



Regular expression examples

- abc

- a|bc

- (a|bb)c

- a[bc]

- a*b

- a?b

- `[^a]*th[aeiou]+[a-z]*`

- `^a`

- `^a.*z$`

- `\bthe\b`

Regular expressions: An extension

- Parentheses allow you to ‘save’ part of a string and access it again.
- In general:
 - `()` creates a group that can be referenced
 - `\n` refers back to the *n*th group in a regular expression
- For example, to search for repeated words: `/([a-z]+) \1 \1/`
- That will find any group of one or more lower case alphabetic characters repeated twice (and separated by a space): ‘I called kitty kitty kitty.’

Regular expressions: An extension

- Note: This extension to Python/Perl/grep/MS “regular expression” syntax actually takes them beyond the realm of regular expressions in the Chomsky hierarchy sense.
- The languages generated by “regular expressions” augmented with this kind of memory device are NOT regular languages, i.e., cannot be recognized by FSAs.
- Why?

So what are regular expressions good for?

- Search
- Search and replace
- Morphotactics: Defining classes of words as sequences of morphemes
- Chunk parsing: Finding units within running text that satisfy patterns such as $D A^* N$

Regular expression substitution in Python

- import re
- re.sub(pattern, replacement, string); use r to mark pattern as a regex
- \g<n> references the nth group

```
>>input = 'The Kiwis is a team.'
>>input =
re.sub(r'(\W)is(\W)', '\g<1>are\g<2>', input)
>>print input
'The Kiwis are a team.'
```

Back to FSAs

- An abstract computing machine
- Consists of a set of states (or nodes in a directed graph) and a set of transitions (labeled arcs in the graph)
- Three kinds of states: plain, start, final

FSA as transition table

- An FSA can also be represented as a transition table:

	Input		
State	a	b	c
0	1	3	-
1:	1	2	3
2:	-	3	-
3:	-	-	-

Recognizing a regular language

FSA's can be used to recognize a regular language.

- Take the FSA and a “tape” with the string to be recognized.
- Start with the start of the tape and the FSA's start state.
- For each symbol on the tape, attempt to take the corresponding transition in the machine.
- If no transition is possible: reject.
- When the string is finished, check whether the current state is a final state.
- Yes: accept. No: reject.

Notes on pseudocode

- Some basic components of algorithms:
 - Loops (e.g., while x is true; for every x in ...)
 - Conditionals (e.g., if ... do this, else do this)
 - Variable assignment (e.g., $i = i + 1$)
 - Evaluating expressions (e.g. if $x = 2$, print $x + 3$)
 - Input values

D-RECOGNIZE in pseudo-code

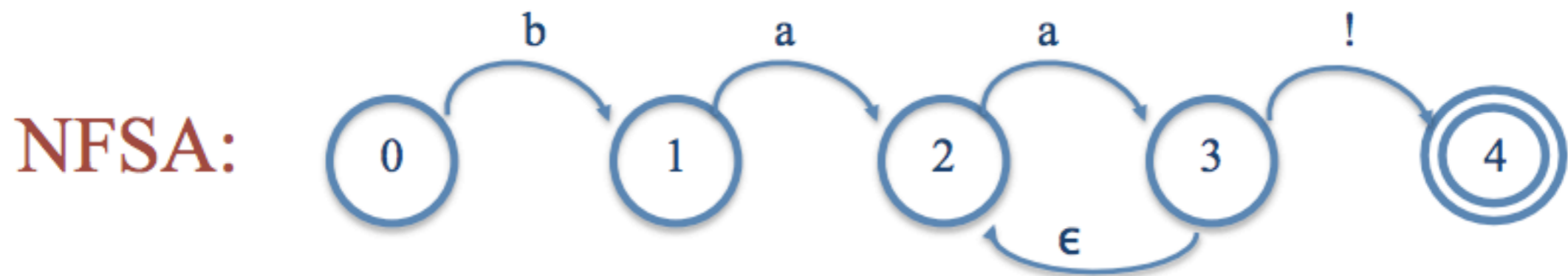
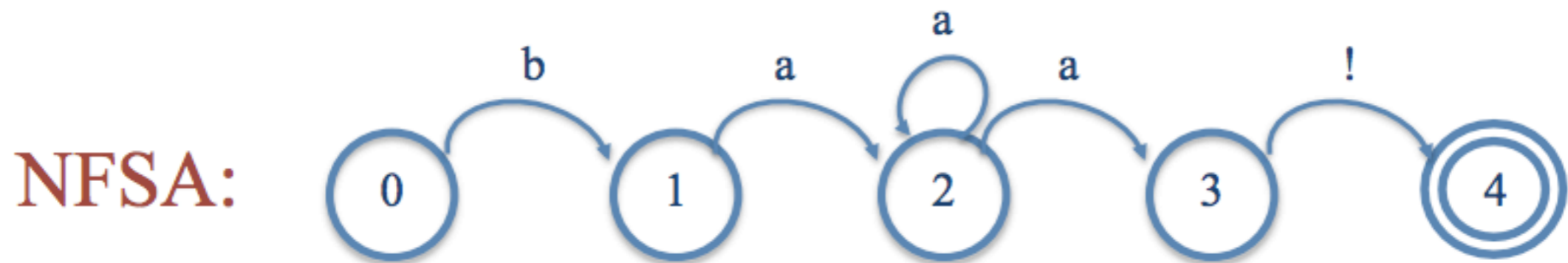
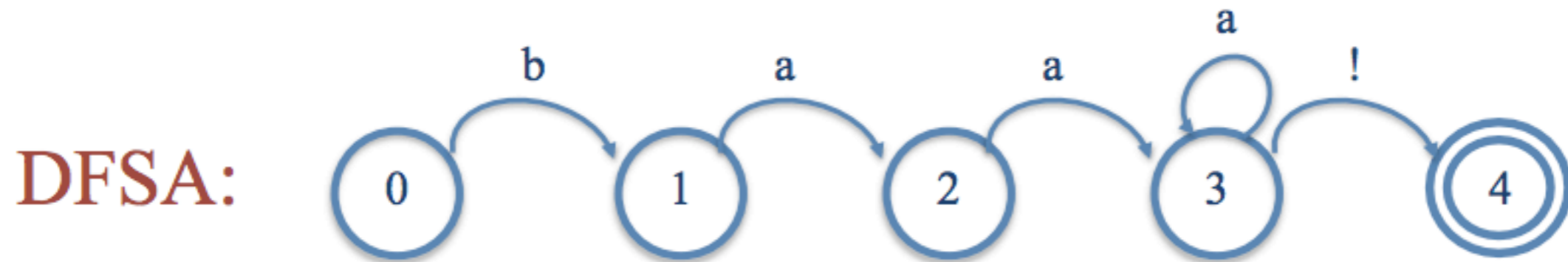
```
function D-Recognize(tape, transition-table) returns accept or reject
  index ← beginning of tape
  current-state ← initial state on transition-table
  loop
    if end of tape has been reached then
      if current-state is an accept state then
        return accept
      else
        return reject
    elseif transition-table[current-state, tape[index]] is empty then
      return reject
    else
      current-state ← transition-table[current-state, tape[index]]
      index ← index + 1
  end
```

Note: This is a slightly modified version of J&M, Figure 2.12, pg29

Non-deterministic FSA (NFSA)

- The FSAs considered so far are deterministic (DFSA): there is only one choice at each node.
- NFSAs include more than one choice at (at least) one node.
- Those choices might include ϵ -transitions, i.e., unlabeled arcs that allow one to jump from one node to another without reading any input.
- Note: Any NFSA can be rewritten as a DFSA.

Equivalent (D)FSA and NFSA



Problem with Non-determinism

- Non-determinism can lead to dead ends when recognizing strings. How?
- Types of solutions:
 - backup
 - look-ahead
 - parallelism
- The backup technique is essentially to record a search-state at each decision point consisting of the position on the input tape and the state (node) of the machine. Then paths to dead ends can be retraced and different decisions made.
- What's the use/point of NFSA's, if there's always an equivalent DFSA?

Overview

- Twitter: tweet classes
- Formal definition of regular languages
- Regular expression examples
- Extension to regular expressions
- Substitution in Python
- Finite state automata
- Recognition with FSA -- pseudocode
- DFSA v. NFSA