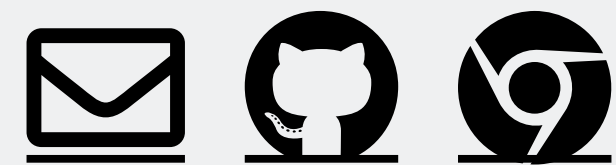

LG 467 Computers in Linguistics

[1-2021] Topic 3: Regular Expressions (and Search)

Sakol Suethanapornkul



Previously...

- Regular expressions (regex)

^ _____

\$ _____

[abc] _____

[^abc] _____

a|b _____

. _____

+ _____

* _____

? _____

\d _____

\w _____

\b _____

\s _____

(abc) _____

a{m} _____

a{m, } _____

a{m,n} _____

Regex in Python

Let's get back to the `re` module. This time, it's `re.split()`

```
source = '''I go on too many dates  
But I can't make 'em stay  
At least that's what people say, mm, mm  
That's what people say, mm, mm'''  
  
print(re.split(r"[', ]", source))
```

Code 5.1

Question: Some things went wrong. What were they?

Regex in Python

You can replace one string with another using `re.sub()`

```
#re.sub(pat, repl, source)

print(re.sub(r"(\w+)?'\w.?", "BAD! ", source))
print(re.sub(r"mm(?!,) ", "urgh", source))
```

Code 5.1

Regex in Python

Lookarounds help you create custom anchors:

Patterns	Matches	Examples
<code>(?!pat)</code>	Anything not followed by pat	<code>r '\w+(?!\d)'</code>
<code>(?<!pat)</code>	Anything not preceded by pat	<code>r '(?<!\d)\w+'</code>
<code>(?=pat)</code>	Whatever that is followed by pat	<code>r '\d+(?=,)'</code>
<code>(?<=pat)</code>	Whatever that is preceded by pat	<code>r '(?<=-)\d+'</code>

nltk_regex_tokenize()

You can use regex to tokenize texts with NLTK:

```
>>> text = 'That U.S.A. poster-print costs $12.40...'
>>> pattern = r'''(?x)          # set flag to allow verbose regexps
...     (?:[A-Z]\.)+          # abbreviations, e.g. U.S.A.
...     | \w+(?:-\w+)*        # words with optional internal hyphens
...     | \$?\d+(?:\.\d+)?%?  # currency and percentages, e.g. $12.40, 82%
...     | \.\.\.             # ellipsis
...     | [][.,;"'()?:_` ]  # these are separate tokens; includes ], [
...     '''
>>> nltk_regex_tokenize(text, pattern)
['That', 'U.S.A.', 'poster-print', 'costs', '$12.40', '...']
```

NLTK Ch.3

See [this](#) for even more complex regex patterns!

Finite-state Automata

Finite-state Automata (FSAs)

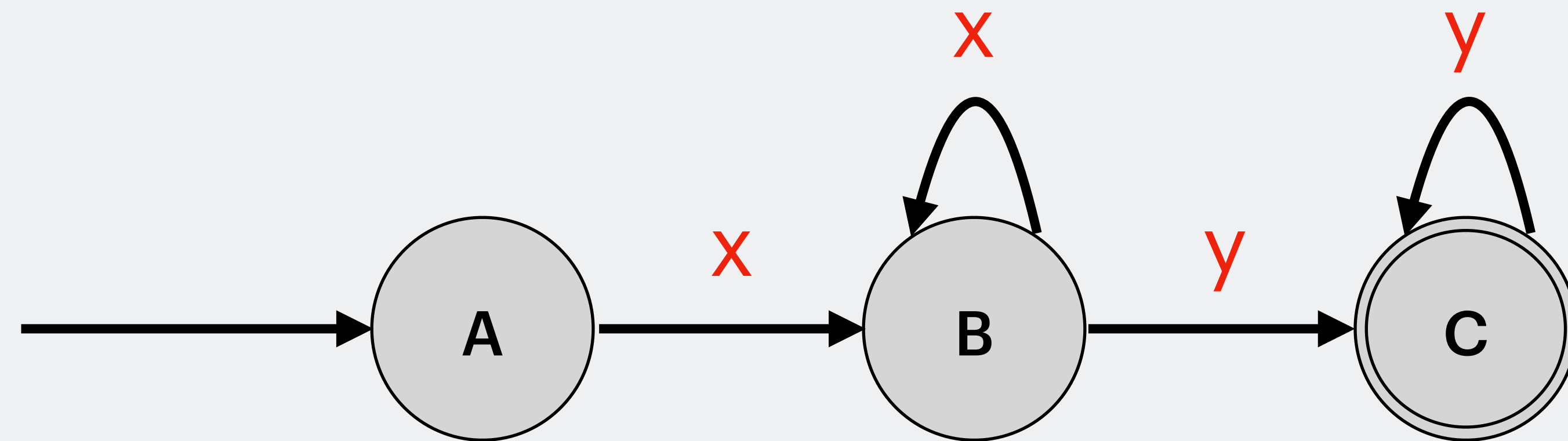
Finite-state automata (FSAs) is a mathematical model:

- that describes one type of language (**regular language**)
- the product of which can be converted to regex (& vice versa)

Finite-state Automata (FSAs)

An automaton comprises four elements:

1. States
2. Initial state
3. Transition rules
4. 1+ final states



xx^*yy^*

Main idea: FSA generates language corresponding to the paths

Finite-state Automata (FSAs)

An automaton can be associated with a set of strings it accepts/generates

- FSAs can be useful tools for recognizing – and generating – subsets of natural language
- But they cannot represent all natural language phenomena

Finite-state Automata (FSAs)

