

# Formal languages, formal grammars, and regular expressions

LING 570

Fei Xia

# Unit #1

- Formal grammar, language and regular expression
- Finite-state automaton (FSA)
- Finite-state transducer (FST)
- Morphological analysis using FST

# Regular expression

- Two concepts:
  - Regular expression in *formal language theory*
  - Regular expression (or **pattern**) in *pattern matching/programming languages*:
    - Ex: `/^\d+(\.\d+)\1/`
- Both concepts describe a set of strings.
- The two concepts are closely related, but the latter is often more expressive than the former.

# Outline

- Formal languages
  - Regular languages
  - Context-free languages
  - ...
- Regular expression in formal language theory
- Formal grammars
  - Regular grammars
  - Context-free grammars
- “Regular expression” in pattern matching

# Formal languages

# Definition of formal language

- An alphabet is a **finite** set of symbols:
  - Ex:  $\Sigma = \{a, b, c\}$
- A string is a **finite** sequence of symbols from a particular alphabet juxtaposed:
  - Ex: the string “baccab”
  - Ex: *empty string*  $\epsilon$
- A formal language is a set of strings defined over some alphabet.
  - Ex1:  $\{aa, bb, cc, aaaa, abba, acca, baab, bbbb, \dots\}$
  - Ex2:  $\{a^n b^n \mid n > 0\}$
  - Ex3: the *empty set*  $\phi$

# Definition of regular languages

- The class of regular languages over an alphabet  $\Sigma$  is formally defined as:
  - The empty set,  $\phi$ , is a regular language
  - $\forall a \in \Sigma \cup \{\epsilon\}$ ,  $\{a\}$  is a regular language.
  - If  $L_1$  and  $L_2$  are regular languages, then so are:
    - (a)  $L_1 \bullet L_2 = \{x y \mid x \in L_1; y \in L_2\}$  (concatenation)
    - (b)  $L_1 \cup L_2$  (union or disjunction)
    - (c)  $L_1^* = \{x_1 x_2 \dots x_n \mid x_i \in L_1, n \in \mathbb{N}\}$  (Kleene closure)
  - There are no other regular languages.

# Kleene star

Another way to define  $L^*$ :

- $L^1 = L$
- $L^n = L^{n-1} \bullet L$
- $L^* = \{ \epsilon \} \cup L^1 \cup L^2 \cup \dots$

Examples:

- $L = \{a, bc\}$
- $L^2 = \{aa, abc, bca, bcbc\}$
- $L^* = \{\epsilon, a, bc, aa, abc, bca, bcbc, aaa, \dots\}$



# Properties

- Regular languages are closed under
  - Concatenation
  - Union
  - Kleene closure
- Regular languages are also closed under:
  - Intersection:  $L_1 \cap L_2$
  - Difference:  $L_1 - L_2$
  - Complementation:  $\Sigma^* - L_1$
  - Reversal

# Are the following languages regular?

- $\{a, aa, aaa, \dots\}$
- Any finite set of strings
- $\{xy \mid x \in \Sigma^*, \text{ and } y \text{ is the reverse of } x\}$
- $\{xx \mid x \in \Sigma^*\}$
- $\{a^n b^n \mid n \in \mathbb{N}\}$
- $\{a^n b^n c^n \mid n \in \mathbb{N}\}$

➔ To prove a language is not regular or context-free, use pumping lemma.

# Regular expression

# Definition of Regular Expression (as in formal language theory)

- The set of regular expressions is defined as follows:
  - (1) Every symbol of  $\Sigma$  is a regular expression
  - (2)  $\epsilon$  is a regular expression
  - (3) If  $r_1$  and  $r_2$  are regular expressions, so are  $(r_1)$ ,  $r_1 r_2$ ,  $r_1 \mid r_2$ ,  $r_1^*$
  - (4) Nothing else is a regular expression.

# Examples

- $ab^*c$
- $a(0|1|2|..|9)^*b$
- $(CV | CCV)^+ C?C?:$  C is a consonant, V is a vowel

Other operations that we can use:

- $a^+ = a a^*$
- $a? = (a | \epsilon)$

# Relation between regular language and Regex

- They are equivalent:
  - With every regular expression we can **associate** a regular language.
  - Conversely, every regular language can be obtained from a regular expression.
- Examples:
  - Regular expression =  $ab^*c$
  - Regular language =  $\{ac, abc, abbc, \dots\}$

# Formal grammars

# Definition of formal grammar

A formal grammar is a concise description of a formal language. It is a  $(N, \Sigma, P, S)$  tuple:

- A finite set  $N$  of *nonterminal symbols*
- A finite set  $\Sigma$  of *terminal symbols* that is disjoint from  $N$
- A finite set  $P$  of *production rules*, each of the form:  
$$(\Sigma \cup N)^* N (\Sigma \cup N)^* \rightarrow (\Sigma \cup N)^*$$
- A distinguished symbol  $S \in N$  that is the *start symbol*



# Chomsky hierarchy

The left-hand side of a rule must contain at least one non-terminal.

$$\alpha, \beta, \gamma \in (N \cup \Sigma)^*, A, B \in N, a \in \Sigma$$

- Type 0: unrestricted grammar: no other constraints.
- Type 1: Context-sensitive grammar:  
The rules must be of the form:  $\alpha A \beta \rightarrow \alpha \gamma \beta$
- Type 2: Context-free grammar (CFGs):  
The rules must be of the form:  $A \rightarrow \alpha$
- Type 3: Regular grammar: The rules are of the forms:  
*right* regular grammar:  $A \rightarrow a, A \rightarrow aB, \text{ or } A \rightarrow \epsilon$   
*left* regular grammar:  $A \rightarrow a, A \rightarrow Ba, \text{ or } A \rightarrow \epsilon$

Are there other kinds of grammars?

# Strings generated from a grammar

- The rules are:  
$$S \rightarrow x \mid y \mid z \mid S + S \mid S - S \mid S * S \mid S / S \mid (S)$$
- What strings can be generated?
- A grammar is **ambiguous** if there exists at least one string which has multiple parse trees.
- Is this grammar ambiguous?

# Languages generated by grammars

- Given a grammar  $G$ ,  $L(G)$  is the set of strings that can be generated from  $G$ .
- Ex:  $G = (N, \Sigma, P, S)$   
 $N = \{S\}, \Sigma = \{a, b, c\}$   
 $P = \{ S \rightarrow a S b, S \rightarrow c \}$

What is  $L(G)$ ?

$$L(G) = \{a^n c b^n\}$$

# The relation between regular grammars and regular languages

- The regular grammars describe exactly all regular languages.
- All the following are equivalent:
  - Regular language: alphabet, operations
  - Regular expression: alphabet, operations
  - Regular grammar: terminals, non-terminals, production rules
  - Finite state automaton (FSA): alphabet, states, edges

# Relation between grammars and languages

Chomsky hierarchy	Grammars	Languages	Minimal automaton
Type-0	Unrestricted	Recursively enumerable	Turing machine
Type-1	Context-sensitive	Context-sensitive	Linear-bounded
Type-2	Context-free	Context-free	Nondeterministic pushdown
Type-3	Regular	Regular	Finite state

# Relation between grammars and languages

(from wikipedia page)\*\*

Chomsky hierarchy	Grammars	Languages	Minimal automaton
Type-0	Unrestricted	Recursively enumerable	Turing machine
n/a	(no common name)	Recursive	Decider
Type-1	Context-sensitive	Context-sensitive	Linear-bounded
n/a	Indexed	Indexed	Nested stack
n/a	Tree-adjoining	Mildly context-sensitive	Thread
Type-2	Context-free	Context-free	Nondeterministic pushdown
n/a	Deterministic context-free	Deterministic context-free	Deterministic pushdown
Type-3	Regular	Regular	Finite state

# How about human languages?

- Are they formal languages?
  - What is the alphabet?
  - What is a string?
- What type of formal languages are they?

crossing dependency:  $N_1 N_2 V_1 V_2$

# Outline

- Formal language
  - Regular language
- Regular expression in formal language theory
- Formal grammar
  - Regular grammar
- **Patterns in pattern matching → J&M-ed2 2.1**



# Patterns in Perl

[ab]	a b
.	match any character
^	the starting position in a string
\$	the ending position in a string
(..)	defines a marked subexpression
a*	match “a” zero or more times
a+	match “a” one or more time
a?	match “a” zero or one time
a{n,m}	“a” appears n to m times

# Special symbols in the patterns

`\s` match any whitespace char

`\d` match any digit

`\w` match any letter or digit

`\S` match any non-whitespace char

...

`\+`, `\-`, `\.`, `\?`, `\*`, ...

# Examples

Integer:  $(\backslash+|\backslash-)?\backslash d+$

Real:  $(\backslash+|\backslash-)?\backslash d+\backslash.\backslash d+$

Scientific notation:  $(\backslash+|\backslash-)? \backslash d+ (\backslash.\backslash d+)?e (\backslash+|\backslash-)?\backslash d+$

Any of the three:

$(\backslash+|\backslash-)? \backslash d+ (\backslash.\backslash d+)? (e (\backslash+|\backslash-)?\backslash d+)?$

# Patterns in Perl and Regex

$/^(.*)\backslash 1\$/ \Leftrightarrow \{xx \mid x \in \Sigma^*\}$

$/^(.+)a(.+)\backslash 1\backslash 2\$/ \Leftrightarrow \{xayxy \mid x, y \in \Sigma^*\}$

➔ The extra power comes from the ability to refer to marked subexpression.