

# Formal Language and Automata Theory (CS21004)

Soumyajit Dey  
CSE, IIT Kharagpur

# Announcements

- The slide is just a short summary
- Follow the discussion and the boardwork
- Solve problems (apart from those we dish out in class)

# Table of Contents

- 1 Pattern Matching
- 2 Regular Expressions
- 3 Regular Grammars

Formal Language  
and Automata  
Theory (CS21004)

Soumyajit Dey  
CSE, IIT  
Kharagpur

Pattern Matching

Regular  
Expressions

Regular Grammars

# Patterns

A *Pattern* captures a family of strings (like FA) : language of the pattern

- $a \in \Sigma : L(a) = \{a\}$  : matches a single symbol
- $\epsilon : \epsilon = \{\epsilon\}$  : matches the null string
- $\phi : L(\phi) = \phi$  : matches empty set  $\phi$
- $\# : L(\#) = \Sigma$  : any symbol
- $@ : L(@) = \Sigma^*$  : any string

Above *atomic* patterns can be operated with unary  $^*$ ,  $^+$ ,  $\neg$  or connected by  $\cup$ ,  $+$ ,  $\cap$ ,  $\circ$  (usually not written, kept silent) to generated other valid patterns

# Patterns

- $x$  matches  $\alpha + \beta$  if  $x$  matches either  $\alpha$  or  $\beta$  :  
 $L(\alpha + \beta) = L(\alpha) \cup L(\beta)$ , similarly you have other rules
- $L(\alpha \cap \beta) = L(\alpha) \cap L(\beta)$ ,  $L(\alpha\beta) = L(\alpha)L(\beta)$ ,  
 $L(\neg\alpha) = \neg L(\alpha) = \Sigma^* - L(\alpha)$ ,
- Can define  $L(\alpha^*)$ ,  $L(\alpha^+)$
- patterns are collections of strings over  
 $\Sigma, \#, @, \epsilon, \phi, \neg, *, +$ ,

Note 1:  $\epsilon, \phi$  and  $\epsilon, \phi$  are different. The boldfaces are symbols in the pattern language

Note 2: '+' is associative over patterns

# Patterns : Applications

- Pattern matching is an important application of FA
- Unix regular exps are basically patterns ( $\Sigma^{??}$ )
- Unix commands 'grep', 'egrep' use FA inside their implementations

# Patterns : examples

- Strings containing at least 3 occurrences of  $a$  :  
 $@a@a@a@$
- all single letters except  $a$  :  $\# \cap \neg a$
- Strings with no occurrences of  $a$  :  $(\# \cap \neg a)^*$

Some books call patterns as regular expressions, we shall make a distinction here.

# Table of Contents

- 1 Pattern Matching
- 2 Regular Expressions
- 3 Regular Grammars

Formal Language  
and Automata  
Theory (CS21004)

Soumyajit Dey  
CSE, IIT  
Kharagpur

Pattern Matching

Regular  
Expressions

Regular Grammars



# Regular Expressions

Represents the idea of patterns as regular language (set) generators ( $\equiv$  FA) in a minimal way using only  $\Sigma, \epsilon, \phi$  as primitive regular expressions and operators  $+, \circ, *$  along with the option of using parenthesis ('(' and ')') to denote operator association.

- $L((aa)^*(bb)^*b) = \{a^{2n}b^{2m+1} \mid n, m \geq 0\}$
- Regex for 'at least one pair of consecutive 0's' :  
 $(0 + 1)^*00(0 + 1)^*$
- Regex for 'no pair of consecutive 0's' :  
 $(1^*011^*)^*(0 + \epsilon) + 1^*(0 + \epsilon) \equiv (1 + 01)^*(0 + \epsilon)$  – how to reason about such equivalence ??

# Regular Expressions, Languages (sets) and FA

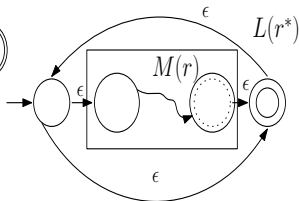
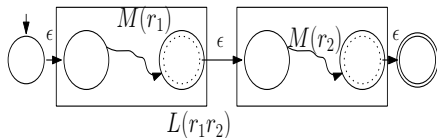
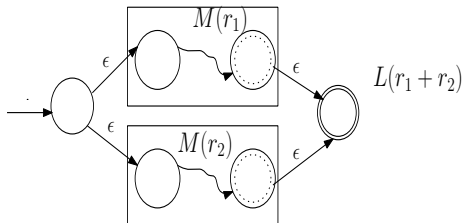
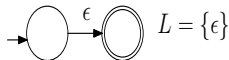
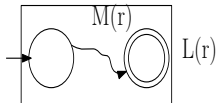
Formal Language  
and Automata  
Theory (CS21004)

Soumyajit Dey  
CSE, IIT  
Kharagpur

Pattern Matching

Regular  
Expressions

Regular Grammars



# Regular Expressions, Languages (sets) and FA

- Automata for primitive regex can be combined as shown
- We can give a formal method for constructing states/transitions of combined machine from states/transitions of simpler machines
- We can prove language equivalence of regex and automata generated as above by induction on number of operators

# DFA to Regex ??

- Let states be  $\{1, 2, \dots, n\}$
- Let  $R(k)_{i,j}$  be the Regex whose language is the set of labels of path from  $i$  to  $j$  without visiting any state with label larger than  $k$ .
- Basis :  $R(0)_{i,j}$  is basically labels of direct paths from  $i$  to  $j$ , i.e.,  $R(0)_{i,j} = a_1 + \dots + a_n$  if  $\delta(i, a_k) = j$  for  $1 \leq k \leq n$ ; Note  $R(0)_{i,i} = \phi$  if there is no self-loop
- Induction :  $R(k)_{i,j} = R(k-1)_{i,j} + R(k-1)_{i,k} \cdot (R(k-1)_{k,k})^* \cdot R(k-1)_{k,j}$

Overall Regex :  $R(n)_{i_0, f_1} + R(n)_{i_0, f_2} + \dots + R(n)_{i_0, f_k}$  with  $i_0$  being the initial state and  $\{f_1, \dots, f_k\}$  being the set of final states.

**COMPLEXITY ???** up to  $n^3$  expressions, each step creates 4 terms for one term.

Using the algorithm,  $\text{NFA} \Rightarrow \text{DFA} \Rightarrow \text{Regex}$  AND  $\text{Regex} \Rightarrow \text{NFA} \Rightarrow \text{DFA}$

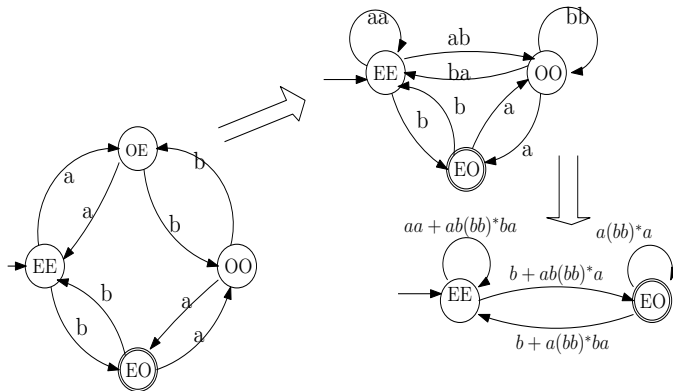
**Observation :** The algorithm we presented can also be devised for NFAs (check Kozen), use  $\Delta$  instead of  $\delta$ .

### Alternate Method :

- Collapse states by allowing regex as transition labels
- Derive regex which connect initial and final state pairs
- Overall expression is the union of such regex

## Example of collapsing

Consider  $\{w \mid n_a(w) \text{ is even}, n_b(w) \text{ is odd}\}$  : difficult to conceive regex directly.



Absence of transition between any state pair  
can be thought of as  $\phi$

Note :  $r\phi = \phi$      $r + \phi = r$      $\phi^* = \lambda$

# Table of Contents

- 1 Pattern Matching
- 2 Regular Expressions
- 3 Regular Grammars

Formal Language  
and Automata  
Theory (CS21004)

Soumyajit Dey  
CSE, IIT  
Kharagpur

Pattern Matching

Regular  
Expressions

Regular Grammars

# Regular Grammars

Regular Language (set)  $\equiv$  FA  $\equiv$  Regular Grammar

- A grammar  $G = (V, \Sigma, S, P)$  is **right linear** if all productions are of the form

$$A \rightarrow xB, A \rightarrow x$$

where  $A, B \in V, x \in \Sigma^*$

- A grammar  $G = (V, \Sigma, S, P)$  is **left linear** if all productions are of the form

$$A \rightarrow Bx, A \rightarrow x$$

where  $A, B \in V, x \in \Sigma^*$

- A regular grammar is one of the two



# Strictly Right Linear Grammars

A grammar  $G = (V, \Sigma, S, P)$  is **strictly right linear** if all productions are of the form

$$A \rightarrow xB, A \rightarrow \lambda$$

where  $A, B \in V, x \in \Sigma \cup \{\lambda\}$

- Any derivation of a word  $w$  from  $S$  has the form

$$S \Rightarrow x_1 A_1 \Rightarrow x_1 x_2 A_2 \Rightarrow \cdots \Rightarrow x_1 \cdots x_n A_n \Rightarrow x_1 \cdots x_n$$

- some  $x_i$  can be  $\lambda$ , connection with NFA ??

♠ For any right-linear grammar  $G$  there exists a strictly right-linear grammar  $H$  such that  $L(G) = L(H)$

# Strictly Right Linear Grammars

If  $G$  is a strictly right-linear grammar, then  $L(G)$  is regular

- Given  $G = (V, \Sigma, P, S)$ , construct NFA  $N = (V, \Sigma, \delta, S, F)$ . The set of states is simply the set of non-terminals of  $G$ . The start state corresponds to the start variable.
- $\delta(A, x) = \{B \mid A \rightarrow xB \in P\}$ ,  $F = \{C \mid C \rightarrow \lambda \in P\}$

To show  $L(G) = L(N)$

- $L(G) \subseteq L(N)$  : by induction on the structure of the derivation
- $L(N) \subseteq L(G)$  : by induction on the structure of the computation

# Strictly Right Linear Grammars

If  $A$  is a regular language, then there is a strictly right-linear grammar  $G$  such that  $L(G) = A$

- If  $L$  is regular,  $\exists$  an NFA  $N = (Q, \Sigma, \delta, q_0, F)$  for  $L$
- We construct a strictly right-linear grammar  $G = (V = Q, \Sigma, P, S = q_0)$  where

$$P = \{A \rightarrow xB \mid B \in \delta(A, x)\} \cup \{C \rightarrow \lambda \mid C \in P\}$$

♠ A language  $A$  is regular **if and only if** there is a strictly right-linear grammar  $G$  such that  $L(G) = A$

# Regular Grammars

Grammar comprising both left-linear and right-linear rules will not necessarily generate a regular language.

Consider

$$S \rightarrow 0A$$

$$S \rightarrow 1B$$

$$S \rightarrow \lambda$$

$$A \rightarrow S0$$

$$B \rightarrow S1$$

The grammar generates  $L = \{ww^R \mid w \in \{0,1\}^*\}$  which is not regular