

Chapter 3: Regular Languages and Regular Grammar

Dr. Yuan Yao

University of Nottingham Ningbo China (UNNC)

Learning Outcomes

Learning outcomes

At the conclusion of this chapter, the students are expected to be able to:

- Identify the language associated with a regular expression.
- Find a regular expression to describe a given language.
- Construct a nondeterministic finite automaton to accept the language denoted by a regular expression.
- Identify whether a particular grammar is regular.
- Construct regular grammars for simple languages.
- Construct an NFA that accepts the language generated by a regular grammar.
- Construct a regular grammar that generates the language accepted by a finite automaton.

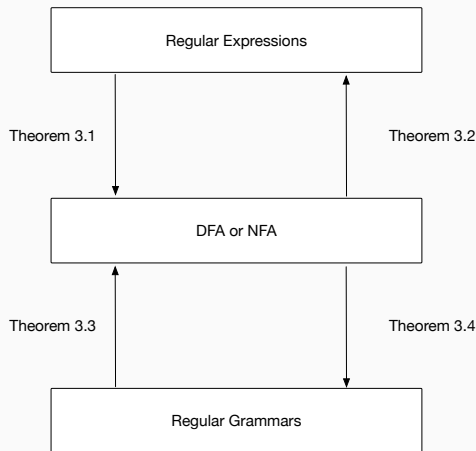
Introduction

Introduction to Regular Languages

- We defined **regular languages** as those that **can be accepted by finite automata**. 正则语言是能被有限自动机 (Finite Automata) 接受的语言。
- In this chapter, we will discuss two other methods for describing regular languages:
 - **Regular expressions.**
 - **Regular grammars.**
- On the **plus** side, regular expressions are **easy to learn** because of their similarity to arithmetic expressions.
- On the **negative** side, there is no obvious way of **extending them to the more complicated classes of languages** we will discuss later.
- Regular grammars, on the other hand, are a special case of more general grammars which we will encounter again in this course.

Introduction to Regular Languages

- FAs, regular expressions, and regular grammars are equivalent.
- Therefore, we can choose whichever method is most appropriate for the situation at hand.
- We will discuss algorithms for converting any of these forms to another later in this chapter.



Regular Expressions

Regular Expressions

- Regular Expressions provide a concise way of describing some languages.
 - Regular Expressions are defined **recursively**. For any **alphabet Σ** :
 - **Primitive regular expressions**:
 - the empty set \emptyset .
 - the empty string λ
 - any symbols $a \in \Sigma$
 - If r_1 and r_2 are regular expressions, then so are:
 - the union: $r_1 \cup r_2$
 - the concatenation: $r_1 r_2$
 - the star-closure: r_1^*
 - parenthesised expression: (r_1)
 - Any string resulting from a **finite** number of these operations on primitive regular expressions is also a regular expression.
- 所有正则表达式必须通过有限次上述操作生成。

Languages Associated with Regular Expressions

- A regular expression r denotes a language $L(r)$.
- Assuming that r_1 and r_2 are regular expressions, then:
 - The regular language \emptyset denotes the empty set.
 - The regular language λ denotes the set $\{\lambda\}$.
 - For any a in the alphabet Σ , the regular expression a denotes the set $\{a\}$.
 - The regular expression $r_1 + r_2$ denotes $L(r_1) \cup L(r_2)$, e.g., $a + b$ means $\{a, b\}$.
 - The regular expression $r_1 r_2$ denotes $L(r_1)L(r_2)$, e.g., ab means $\{ab\}$.
 - The regular expression r_1^* denotes $(L(r_1))^*$, e.g., a^* means $\{\lambda, a, aa, \dots\}$.
 - The regular expression (r_1) denotes $L(r_1)$.

运算符优先级：括号 > 星闭包 > 连接 > 并集（如 $a + b^*$ 等价于 $a + (b^*)$ ）。

Exercise: Regular Expression

- What is the language $L(ab^* + c)$?

Determining the Language Denoted by a Regular Expression

- By combining regular expressions using the given rules, arbitrarily complex expressions can be constructed.
- In applying operations, we observe the following **precedence rules**:
 - star closure precedes concatenation:
 - Example: ab^* should be interpreted as $a(b)^*$ rather than $(ab)^*$.
 - Thus, $L(ab^*) = \{a, ab, abb, \dots\}$.
 - concatenation precedes **union**:
 - Example: $ab + c$ should be interpreted as $(ab) + c$ rather than $a(b + c)$.
 - Thus, $L(ab^* + c) = \{c, a, ab, abb, \dots\}$.
 - Parentheses are used to override the normal precedence of operators.
 - Hence, the language $\{ab, ac\}$ is generated by the regular expression $a(b + c)$.

Sample Regular Expressions and Associated Languages

- $(ab)^*$
- $a + b$
- $(a + b)^*$
- $a(bb)^*$
- $a^*(a + b)$
- $(aa)^*(bb)^*b$
- $(0 + 1)^*00(0 + 1)^*$

Sample Regular Expressions and Associated Languages

- $L((ab)^*) = \{(ab)^n | n \geq 0\}$
- $L(a + b) = \{a, b\}$
- $L((a + b)^*) = \{a, b\}^*$ or any strings formed with a and b
- $L(a(bb)^*) = \{ab^{2n} | n \geq 0\}$
- $L(a^*(a + b)) = \{a^n a, a^n b | n \geq 0\}$
- $L((aa)^*(bb)^*b) = \{a^{2n}b^{2m}b | n \geq 0, m \geq 0\}$
- $L((0 + 1)^*00(0 + 1)^*)$: Binary strings containing at least one pair of consecutive zeros. 包含至少连续两个0的二进制串

Equivalence of Regular Expressions

- Two regular expressions are equivalent if they denote the same language.
- For example, $(a + b)^*$ and $(a^*b^*)^*$ are equivalent, because:

- $L((a + b)^*) = L((a^*b^*)^*) = \{a, b\}^*$

- Another interesting case, what are the languages for the following regular expression?

- $r_1 = (1^*011^*)^*(0 + \lambda) + 1^*(0 + \lambda)$

- $r_2 = (1 + 01)^*(0 + \lambda)$

o $(1^*011^*)^*$: 任意数量的"含单个0的1串"模式

o $(1 + 01)^*$: 重复选择1或01

o $(0 + \lambda)$: 可选结尾0

o $(0 + \lambda)$: 可选结尾0

o $1^*(0 + \lambda)$: 纯1串可选结尾0

Equivalence of Regular Expressions

- Two regular expressions are equivalent if they denote the same language.
- For example, $(a + b)^*$ and $(a^*b^*)^*$ are equivalent, because:
 - $L((a + b)^*) = L((a^*b^*)^*) = \{a, b\}^*$
- Another interesting case, what are the languages for the following regular expression?
 - $r_1 = (1^*011^*)^*(0 + \lambda) + 1^*(0 + \lambda)$
 - $r_2 = (1 + 01)^*(0 + \lambda)$
 - $L(r_1) = L(r_2) = \{w \in \{0, 1\}^* \mid w \text{ has no pair of consecutive zeros} \}$
- We have seen something similar...

Equivalence of Regular Expressions

- Two regular expressions are equivalent if they denote the same language.
- For example, $(a + b)^*$ and $(a^*b^*)^*$ are equivalent, because:
 - $L((a + b)^*) = L((a^*b^*)^*) = \{a, b\}^*$
- Another interesting case, what are the languages for the following regular expression?
 - $r_1 = (1^*011^*)^*(0 + \lambda) + 1^*(0 + \lambda)$
 - $r_2 = (1 + 01)^*(0 + \lambda)$
 - $L(r_1) = L(r_2) = \{w \in \{0, 1\}^* \mid w \text{ has no pair of consecutive zeros} \}$
- The complement of r_1 or r_2 , $r = (0 + 1)^*00(0 + 1)^*$
- There is no obvious connection between r and $r_1, r_2...$

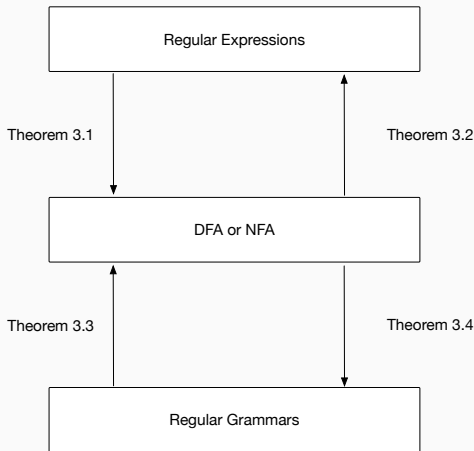
- One of the most widely accessible applications of regular expressions is in search and pattern matching. 搜索与模式匹配
- Any non-trivial editor provides search using regular expressions, e.g., Emacs, Vim...

Connections Between Regular Expression and Regular Language

Regular Expression and NFA

对任意正则表达式 r ，存在接受 $L(r)$ 的NFA

- **Theorem 3.1:** For any regular expression r , there is a nondeterministic finite automaton that accepts the language denoted by r .
- Since NFA and DFA are equivalent, for any regular expression r , the language $L(r)$ is also regular.
- How to prove this theorem?



A Constructive Proof

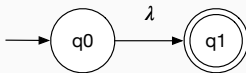
- Construction of an NFA to accept a language $L(r)$ where r is a regular expression.
- We start with **primitive regular expressions**:
- Draw a NFA with two states for the following regular expressions:
 - the empty set.
 - the empty string.
 - Any individual symbol $a \in \Sigma$.

A Constructive Proof

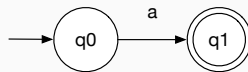
- Construction of an NFA to accept a language $L(r)$ where r is a regular expression.
- We start with **primitive regular expressions**:
- Draw a NFA with two states for the following regular expressions:
 - (a) the empty set.
 - (b) the empty string.
 - (c) any individual symbol $a \in \Sigma$.



(a)



(b)



(c)

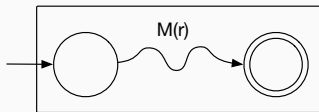
- Before going forward, we need to prove the following claim:
- **Claim:** for every NFA with arbitrary number of final states, there is an equivalent NFA with only one final state.
对于任意终态数量的NFA，存在仅含一个终态的等价NFA

A Constructive Proof

- Before going forward, we need to prove the following claim:
- **Claim:** for every NFA with arbitrary number of final states, there is an equivalent NFA with only one final state.
- **Hint:** Introduce a new final state p_f . For every state $q \in F$, add a λ -transition from q to p_f , i.e., $\delta(q, \lambda) = \{p_f\}$. Make p_f the only final state. Then prove that $\delta^*(q_0, w) \in F$ in the original NFA if and only if $\delta^*(q_0, w) = \{p_f\}$ after the modification. So, these two NFAs are equivalent.
- How about the same claim but for DFA? is it true?

A Constructive Proof

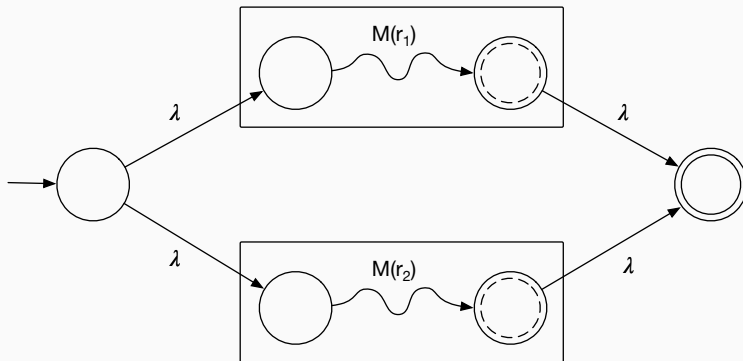
- We could therefore use the following representation for automata $M(r)$ that accept the language $L(r)$ denoted by a regular expression r .



- Assume r_1 and r_2 are two regular expressions, then how to construct an automaton that accept $L(r_1 + r_2)$?

A Constructive Proof

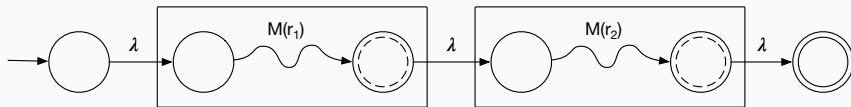
- Assume r_1 and r_2 are two regular expressions, then how to construct an automaton that accept $L(r_1 + r_2)$?



- Assume r_1 and r_2 are two regular expressions, then how to construct an automaton that accept $L(r_1r_2)$, i.e., the concatenation?

A Constructive Proof

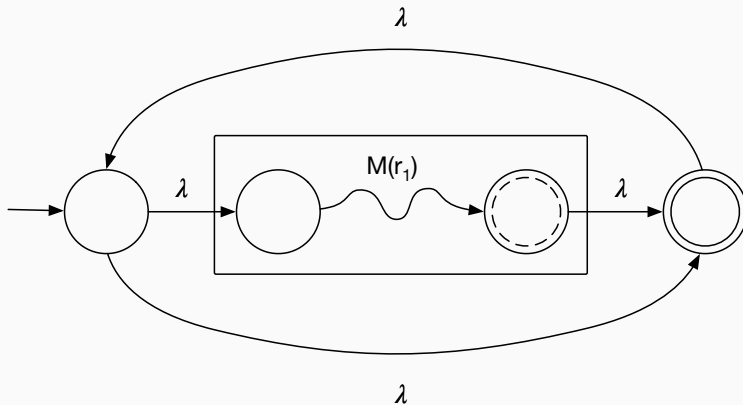
- Assume r_1 and r_2 are two regular expressions, then how to construct an automaton that accept $L(r_1r_2)$, i.e., the concatenation?



- Assume r_1 is a regular expressions, then how to construct an automaton that accept $L(r_1^*)$?

A Constructive Proof

- Assume r_1 is a regular expressions, then how to construct an automaton that accept $L(r_1^*)$?



Exercise: Regular Expression to NFA

- Given the following regular expression

$$r = (a + bb)^*(ba^* + \lambda)$$

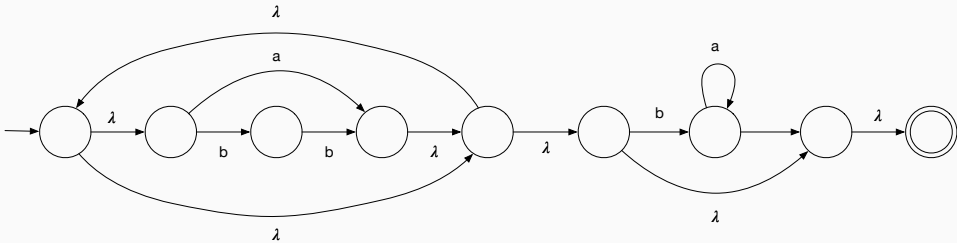
- Construct a DFA that accept $L(r)$.

Exercise: Regular Expression to NFA

- Given the following regular expression

$$r = (a + bb)^*(ba * + \lambda)$$

- Construct a DFA that accept $L(r)$.



Regular Expressions for Regular Languages

- So far, we have learnt that, for any regular expression r , the language $L(r)$ is a regular language.
- What about the other direction?
- **Question:** Is it true that for any regular language L , there exists a regular expression r such that $L = L(r)$?
- **Theorem 3.2** Let L be a regular language. Then there exists a regular expression r such that $L = L(r)$.
- How to prove it?

Regular Expressions for Regular Languages

- So far, we have learnt that, for any regular expression r , the language $L(r)$ is a regular language.
- What about the other direction?
- **Question:** Is it true that for any regular language L , there exists a regular expression r such that $L = L(r)$?
- **Theorem 3.2** Let L be a regular language. Then there exists a regular expression r such that $L = L(r)$.
- How to prove it?
- To prove theorem 3.2, we need an algorithm which, given a finite automaton M , returns a regular expression r such that $L(r) = L(M)$.
- We are not going to introduce the details, please refer to the textbook. This algorithm is also implemented in JFLAP.

Regular Grammar

Regular Grammar

- The production rule for the following two languages are as follows:

- $L_1 = \{a^n b^n | n \geq 0\}$

$$S \rightarrow aSb \mid \lambda$$

- $L_2 = \{a^n b^m | n, m \geq 0\}$

$$S \rightarrow aS \mid A$$

$$A \rightarrow bA \mid \lambda$$

- **Question:** Can we identify the type of grammars that generates **regular language**?
- What's the difference between the production rules above?

Regular Grammar

- The production rule for the following two languages are as follows:

- $L_1 = \{a^n b^n | n \geq 0\}$

$$S \rightarrow aSb \quad | \quad \lambda$$

- $L_2 = \{a^n b^m | n, m \geq 0\}$

$$S \rightarrow aS \quad | \quad A$$

$$A \rightarrow bA \quad | \quad \lambda$$

- What's the difference between the production rules above?
 - in the production $S \rightarrow aSb$, symbols appear on both sides of S .
 - whereas in the production rules of the second grammar, symbols are added only to one side (in this case, the left side) of the variable.

Right-Linear and Left-Linear Grammar

- **Right-linear Grammar:** a grammar $G = (V, T, S, P)$ is said to be right-linear if all productions in P are of the form:

$$A \rightarrow xB$$

$$A \rightarrow x$$

where $A, B \in V$ and $x \in T^*$

- **Left-linear Grammar:** a grammar $G = (V, T, S, P)$ is said to be left-linear if all productions in P are of the form:

$$A \rightarrow Bx$$

$$A \rightarrow x$$

where $A, B \in V$ and $x \in T^*$

- We say that a grammar is regular if it is either right-linear or left-linear.

Exercise: Regular Grammar

- Given the following grammar G :

$$G = (\{S\}, \{a, b\}, S, P)$$

where P is defined as follows:

$$S \rightarrow abS \mid a$$

- Is the above grammar left-linear or right-linear or neither?
- Can you give a regular expression for the regular language generated by the above grammar?

Exercise: Regular Grammar

- Given the following grammar G :

$$G = (\{S\}, \{a, b\}, S, P)$$

where P is defined as follows:

$$S \rightarrow abS \mid a$$

- Is the above grammar left-linear or right-linear or neither?
 - **Right-linear.**
- Can you give a regular expression for the regular language generated by the above grammar?
 - **$(ab)^*a$**

Right-Linear Grammars Generate Regular Languages

给定右线性文法 $G=(V,T,S,P)$, 其生成的语言 $L(G)$ 是正则语言。

- Theorem 3.3 Let $G = (V, T, SP)$ be a right-linear grammar. Then $L(G)$ is a regular language.
- We could prove this theorem constructively.
- In fact, there is an algorithm for constructing an NFA to accept the language generated by a given right-linear grammar G .

Right-Linear Grammars Generate Regular Languages

- How to construct an NFA to accept the language generated by a given right-linear grammar G :
 - Label the NFA start state with S and a final state V_f .
 - For every variable symbols V_i , create an NFA state and label it V_i .
 - For each production of the form $A \rightarrow aB$, label a transition from state A to B with symbol a .
 - For each production of the form $A \rightarrow a$, label a transition from state A to V_f with symbol a .
 - Note: you need to add intermediate states for productions with more than one terminal on the right-hand side.
- **Question:** Can you see why in general we get an NFA, rather than a DFA?

Exercise: Construct an NFA to accept the given regular grammar

- Given the regular grammar $G = (\{V_0, V_1\}, \{a, b\}, V_0, P)$, where P is defined as follows:

$$V_0 \rightarrow aV_1$$

$$V_1 \rightarrow abV_0 \mid b$$

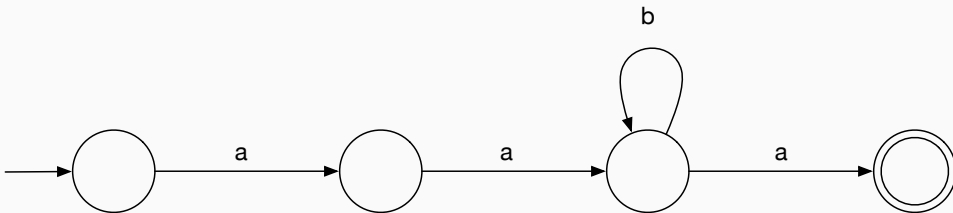
Right-Linear Grammars for Regular Languages

对于字母表 Σ 上的任意正则语言 L , 存在右线性文法 $G=(V, T, S, P)$ 满足 $L = L(G)$

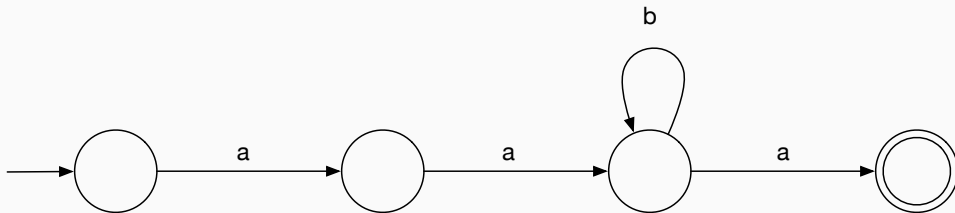
- Theorem 3.4: If L is a regular language on the alphabet Σ , then there exists a right-linear grammar $G = (V, T, S, P)$ such that $L = L(G)$.
- How to prove it?
 - There is an algorithm that, given any DFA M accepting a regular language L , constructs a right-linear grammar G which generates the same language:
 - Each state in the DFA corresponds to a variable symbol in G .
 - For each DFA transition from state A to state B labeled with symbol a , there is a production of the form $A \rightarrow aB$ in G .
 - For each final state F_i in the DFA, there is a corresponding production $F_i \rightarrow \lambda$ in G .

Exercise

- Given the following finite automaton M , write a regular grammar G such that $L(M) = L(G)$.



Exercise



- The grammar $G = (\{A, B, C, D\}, \{a, b\}, A, P)$, where P is defined as follows:

$A \rightarrow aB$

$B \rightarrow aC$

$C \rightarrow bC \mid aD$

$D \rightarrow \lambda$

Notes: Regular Grammars

- Remember that a grammar is said to be regular if it is either left-linear or right-linear.
- For simplicity, we stated all our results in terms of right-linear grammars, but very similar arguments prove the corresponding results for left-linear grammars.
- Caution:** A grammar is regular if either all its productions are left-linear, or they are all right-linear, but not a combination of both.
- For example, consider the grammar $G = (\{S, A, B\}, \{a, b\}, S, P)$ with productions:
 - $S \rightarrow A$
 - $A \rightarrow aB \mid \lambda$
 - $B \rightarrow Ab$
- Question:** What is the language $L(G)$?

1. 产生式展开:

$S \rightarrow A \rightarrow aB \rightarrow aAb \rightarrow aaBb \rightarrow aabAb \rightarrow \dots$

无法形成有限推导，产生无限嵌套模式

2. 语言特性:

- 无法生成任何终结符字符串
- 所有推导都会陷入非终结符循环
- 最终结论: $L(G) = \emptyset$ (空语言)

Application: Text Editing and Pattern Matching

- We have already mentioned that, almost any sophisticated text editor (such as Emacs, Vim, Netbeans, etc.) allows search by regular expressions.
- **Question:** How is it possible to perform search using a regular expression as input?
- The text editor goes through the following steps:
 1. Convert the regular expression into an equivalent NFA (Theorem 3.1).
 2. Convert the NFA to an equivalent DFA.
 3. Minimize the DFA (which we do not discuss in this module).
 4. Finally, run the DFA over the input for pattern matching.
 1. 将该正则表达式转换为等效的 NFA (定理 3.1)。
 2. 将非确定有限自动机转换为等效的确定有限自动机。
 3. 简化确定有限自动机 (在本模块中我们不对其进行详细讨论)。
 4. 最后, 对输入数据运行确定有限自动机 (DFA), 以进行模式匹配。

Application: Compilation and Pattern Matching

- **Question:** How is it possible for a C compiler to check whether a string of symbols is a valid identifier?
- The designers of the compiler must go through the following steps:
 1. Convert the regular grammar to an equivalent NFA (Theorem 3.3).
 2. Convert the NFA to an equivalent DFA.
 3. Minimize the DFA (which we do not discuss in this module).
 4. Finally, incorporate the DFA into the compiler for pattern matching.
- **As can be seen, the relatively simple results that we have discussed so far are used in practice for crucial applications (e.g., pattern matching in text editors and compilers).**