# 西安财经学院 信息学院

<u>　　编译原理　　</u>　实验报告

姓名：陈伯硕
学号：1831050010
班级：计本 1801
指导教师：李薇

成绩：_____

**实验名称** 正规式、自动机的相互转换

**实验日期：** 2020 年 10 月 26 日

## 一、实验目的

1. 理解正规式、NFA、DFA、DFA 最小化的基本原理；
2. 掌握正规式向 NFA 转换的算法；
3. 掌握 NFA 向 DFA 转换的算法；
4. 掌握 DFA 最小化的算法。

## 二、实验内容

1. 设计并实现正规式向 NFA 转换的算法；
2. 设计并实现 NFA 向 DFA 转换的算法；
3. 设计并实现 DFA 最小化的算法。

## 三、实验要求

1. 从下列实验内容中任选一个或几个算法实现；
   （1）设计并实现正规式向 NFA 转换的算法：输入为正规式，输出为一个 NFA。
   （2）设计并实现 NFA 向 DFA 转换的算法：输入为一个 NFA 五元组，输出一个与其等价的 DFA 五元组。
   （3）设计并实现 DFA 最小化的算法：输入为任意一个 DFA，输出此 DFA 等价的状态最少的 DFA。
2. 要求根据算法要求设计合理的数据结构；
3. 编程实现转换过程。

# 四、 算法设计

**1. 自动机的数据结构与算法**

**(1) 自动机的数据结构**

自动机由元组 $(\Sigma, S, S_0, F, f)$ 组成，其中

1. $\Sigma$ 为字符的集合, 空串 $\epsilon \notin \Sigma$
2. $S$ 为状态集合
3. $S_0 \in S$ 为初始态
4. $F \subset S$ 是终态的集合
5. $f : S \times (\Sigma \cup \{\epsilon\}) \to S$ 为状态转换函数, 在程序中可以使用哈希表 $\text{hash}(S_i, S_{j_k}) = C_{ij_k}$, $\forall c_{ij_k} \in C_{ij_k}$, 使得 $f(S_i, c_{ij_k}) = S_{j_k}$ 在 python 中，hash 函数由字典数 *dict* 据类型直接实现, 对于两个变量的哈希表，可以用字典嵌套字典实现

**(2) 自动机的算法**

设正则表达式 $s, t$，其对应 NFA 为 $N(s), N(t)$

对于表达式 $s|t$ 对于运算链接 (concatenation)$N(s)N(t)$, 对应的算法如算法 1

---

**算法 1** concatenation 运算

---

concatenation($N(s)$,$N(t)$)

1   renameState($N(t)$) $\triangleright$ 修改 $N(t)$ 状态名，防止歧义
2   merge($N(s), N(t)$) $\triangleright$ 将 $N(t)$ 的状态和转换函数复制给 $N(s)$
3   $F_s \leftarrow F_t \triangleright N(s)$ 的终态变为 $N(t)$ 的终态 $F_t$
4   free($N(t)$) $\triangleright$ 释放 $N(t)$ 的空间
5   **return** $N(s)$

---

并运算 (union)$N(s)|N(t)$

union($N(s), N(t)$)

1   rename($N(s)$)

2   rename($N(t)$)

3   copy($N(s), N(t)$)

4   **new** $S_0 \triangleright$ 新初态

5   **new** $S_t \triangleright$ 新终态

6   addEdge($S_0, S_0^{(s)}, \epsilon$) $\triangleright$ $S_0^{(s)}$ 为原 $N(s)$ 的初态

7   addEdge($S_0, S_0^{(t)}, \epsilon$) $\triangleright$ $S_0^{(s)}$ 为原 $N(t)$ 的初态

8   **for** $S_f \in F^{(s)} \cap F^{(t)}$:

9       addEdge($S_f, S_t, \epsilon$)

10  $F \leftarrow \{S_t\} \triangleright$ 新的终态集

11  free($N(t)$)

12  **return** $N(s)$

对于 $s*$ 的形式，表示 $s$ 匹配一次或多次，对自动机来说，无条件地进入终态，无条件返回初态，如算法 3

**算法 3** star 运算

star($N(t)$)

1  **for** $S_f \in F$:

2      addEdge($S_f, S_0, \epsilon$)

3      addEdge($S_0, S_f, \epsilon$)

4  **return** $N(t)$

## 2. 正则表达式解析

正则表达式的结构可以由以下解析式组成

```
<regex>  ::= <term> '|' <regex>
         | <term>


<term>   ::= { <factor> }


<factor> ::= <base> { '*' }
```

```
<base>   ::= <char>
         |  '\' <char>
         |  '(' <regex> ')'
```

这些算法可以用递归实现，我们首先完成以下函数

1. peek() 查看模式串中下一个字符
2. eat(*c*) 检查下一个字符是否为 *c*，如果是，将模式串中去掉该元素，即返回第一个字符之后的元素
3. next() 返回下一个元素，并在模式串中去掉该元素

我们从最简单的部分开始，首先通过算法 4 检测 base, 处理转义符，括号和普通单个字母

---

**算法 4** base 部分的解析

---

base()

1  **if** peek() = '(':        ▷ 匹配括号
2      eat('(')                ▷ 处理掉'('
3      r ← regex()            ▷ 这一部分为另一个正则表达式, 在算法 7实现
4      eat(')')
5      **return** r            ▷ 此时括号内的表达式的 NFA 是解析结果
6  **else if** peek() = '\':    ▷ 处理转义符之后的字符
7      eat( '\' )
8      esc ← next()           ▷ 获得转义符之后的一个字符
9      **return** basicConstruct(esc)    ▷ 创建只含有 *esc* 的 NFA
10  **else** :
11      **return** basicConstruct(next())    ▷ 只有一种字母的情况

---

factor 部分由 base 和若干个 * 组成，我们对算法 4 的结果进行 star 运算 (算法 3) 具体过程如算法 5

**算法 5** factor 部分的解析

---

factor()

s 1    base ← base()   ▷ 从算法 4 中得到 * 之前部分的 NFA

  2    **while** parttern **and** peek() = '*':   ▷ 处理所有'*' 字符，并进行相应运算

  3       eat('*')

  4       star(base)   ▷ 对 base 进行 star 运算 (算法 3)

  5    **return** base

---

term 部分由若干个 factor 组成，他们之间用 concatenation 运算链接 (算法 1) 具体实现如 算法 6

**算法 6** term 部分的解析

---

term()

1    term ← basicConstruct($\epsilon$)   ▷ 只有空串转换的 NFA

2    **while** pattern **and** peek() ≠ ')' **and** peek() ≠ '*': ▷ 下一个字符不能是需要运算符号

3       f ← factor() ▷ 读取下一个 term(算法 5)

4       term ← concatenation(term, f) ▷ 更新 term

5    **return** term

---

最终，我们可以将一个正则表达式分解为 term 或者 term '|' regex 的形式我们对相应 NFA union 运算 (算法 2) 然后得到最终的 NFA 具体细节见算法 7

**算法 7** regex 部分的解析

---

regex()

1    term ← term() ▷ 第一部分为 term 部分

2    **if** pattern **and** peek() = '|':

3       regex ← regex() ▷ 下一部分是另一个正则表达式，用自身解析

4    **else** :

5       **return** term

---

# 五、 实验步骤与结果

## 1. 代码测试自动机的基本构造

构造自动机类的数据类型和相关方法，并检查状态转换是否正确，运行

```
test = Automata('ab')
test.set_start_state(1)
test.add_final_states(2)
test.add_final_states(2)
test.add_transition(1,2,set(['a','b']))
test.add_transition(1,3,set('b'))
test.draw('../docs/figures/test_automata.pdf')
```

其中打印显示数据类型如下

```
1  {1: {2: {'a', 'b'}, 3: {'b'}}}
2  states: {1, 2, 3}
3  start state:  1
4  final state:  {2}
5  transitions:
6  1->2 on 'a'
7  1->2 on 'b'
8  1->3 on 'b'
```
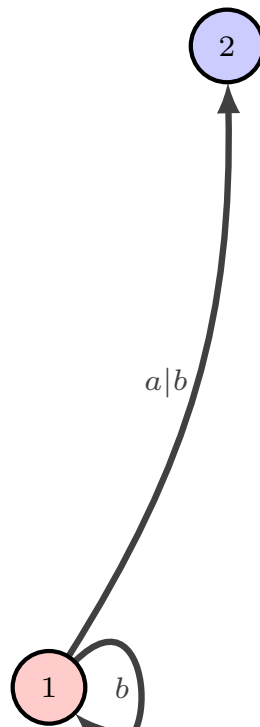
绘图结果如图 1



**图 1 测试一个自动机的数据结构**

图 1蓝色点代表终态，红色代表初态，其他状态为绿色

## 2. 测试自动机的方法和运算

接下来测试自动机的运算，对于图 2a的自动机 $N_1$



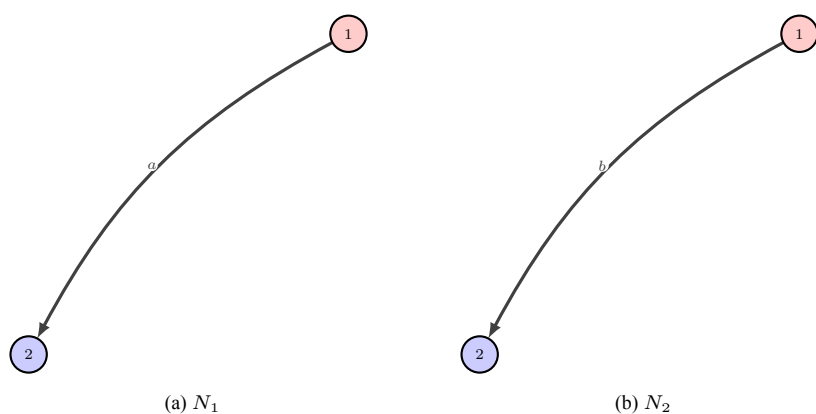(a) $N_1$                    (b) $N_2$

**图 2　简单的自动机**

根据算法 3 测试 $\text{star}(N(s))$ 结果如



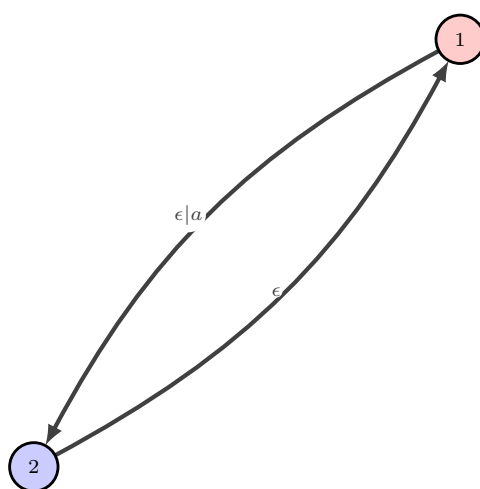**图 3　star 运算测试结果**

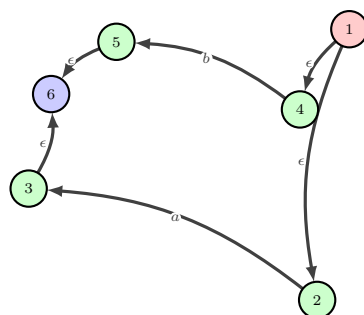接下来完成算法 2 并测试 $N(s), N(t)$ 如图 4

**图 4　union 运算测试**

最后根据算法 1 完成 concatenation 运算代码并绘图，得



**图 5　concatenation$(N(s), N(t))$ 的结果**

## 3. 测试正则表达式解析

根据 算法 (4) ~ (7) 完成正则表达式解析算法我们测试正则表达式 $a|b$ 得到图 6

**图 6    正则表达式 $a|b$ 生成的 NFA**

图 6 中 $6 \to 7,8 \to 11$ 展示了状态转换的逻辑，但是算法中引入大量空串，使生成 $NFA$ 变得复杂，对于像 $(a|b)*ab$ 这样复杂的表达式 () 视觉上很难识别转换关系需要进一步化简。



**图 7    比较复杂的表达式 $(a|b)*ab$**

# 六、 实验总结

这部分实验是一组算法，非常锻炼逻辑能力和设计能力，对数据类型的封装可以应用许多面相对象技术，如实现运算时，可以使用一个对象的方法更新，也可以使用静态方法，基本的自动机类型可以由工厂模式的类方法构造，这一些对以后的程序设计很有帮助。

算法方面，文字的分析很好的展现了递归的简单优雅，利用递归可以快速构造代码，且符合思考的逻辑，但是如果想进一步优化算法，需要利用栈来做进一步分析。

这一类算法的调试，分析都是从小的部分自顶而上构造，自顶而下的思考分析，锻炼了逻辑思维。

对于网上的算法，需要自己的理解分析对于网上的代码需要改造为自己理解的形式，直接不假思索的复制，沉迷于低质量的信息，在我的实验过程中看到了一些经典的讲解 [2] 和一些完整的实现 [3] 这一些都是参考学习的优质资源，比许多明显是别的学生作业的资源优质许多，这些可以对照我们学过的理论 [1] 对编译系统有更深的了解，进而熟悉一套解析的逻辑，增强自己的能力

# 参考文献

[1] Alfred V Aho, Ravi Sethi, and Jeffrey D Ullman. *Compilers, principles, techniques*, volume 7. Addison Wesley, 2 edition, 2006.

[2] mattmight. Parsing regular expressions with recursive descent. `http://matt.might.net/articles/parsing-regex-with-recursive-descent/`.

[3] sdht0. automata-from-regex. `https://github.com/sdht0/automata-from-regex`.

# 附录 A　依赖的安装

本文代码的 Atutomata.draw() 方法依赖

```
network2tikz
```

可以通过以下命令在清华镜像网站安装全部

```
pip install -i http://mirrors.aliyun.com/pypi/simple \
--trusted-host mirrors.aliyun.com/pypi/simple/ \
-r requirements.txt
```

如果在程序中不使用画图功能，可以忽略安装

# 附录 B   自动机的相关方法与测试源代码

```python
"""
filename src/Automata.py
reference https://github.com/sdht0/automata-from-regex/blob/master/AutomataTheory.py
"""
from __future__ import annotations # type hint within a class
from typing import *
# see https://stackoverflow.com/questions/41135033/type-hinting-within-a-class
from matplotlib import pyplot as plt
import networkx as nx
import matplotlib as mpl
plt.rcParams.update({
    "text.usetex": True,
    "font.family": "sans-serif",
    "font.sans-serif": ["Helvetica"]})
# for Palatino and other serif fonts use:
plt.rcParams.update({
    "text.usetex": True,
    "font.family": "serif",
    "font.serif": ["Palatino"],
})


class Automata:
    """
    class to represent a automata

    :param input_alphabet: a set of input symbols
    :type input_alphabet: set,optional

    :ivar empty_string: empty string, denoted by :math:`\epsilon`
    :ivar self.states: a finite states of S
    :ivar self.transitions: 11
    :ivar self.input_alphabet: a set of input symbols
    :ivar self.final_states: the set of final state
    :ivar self.transitions: the transitions functions,
        `translations[f][t] = d` where f is from state,t in to state,
        d is the dict of states where d[state] = set of input symbols
    """
    empty_string = set([r'\epsilon'])

    def __init__(self, input_alphabet: set):
        self.states = set() # a finite states of S
        self.input_alphabet = input_alphabet # a set of input symbols
        self.start_state = None
```

```python
45        self.final_states = set()
46        self.transitions = dict()
47
48    # @staticmethod
49    # def empty_string() -> str:
50    #     r"""get the symbol of empty_string symbol :math:`\epsilon`
51
52    #     :return: r'\epsilon'
53    #     :rtype: str
54    #     """
55    #     return r'\epsilon'
56
57    def set_start_state(self, state: int):
58        """set the start state
59
60        :param state: the label of start state
61        :type state: int
62        """
63        self.start_state = state
64        self.states.add(state)
65
66    def add_final_states(self, *states):
67        """add the final states
68
69        :param states: the list of states
70        """
71        for state in states:
72            self.final_states.add(state)
73
74    def add_transition(self, from_state: int, to_state: int, input_symbols: set):
75        """add the transition to transfer functions
76        (`self.transitions` in the program)
77
78        :param from_state: the begin state
79        :type from_state: int
80        :param to_state: the next state
81        :type to_state: int
82        :param input_symbols: the transfer symbols to the next states
83        :type input_symbols: set
84        """
85
86        self.states.add(from_state)
87        self.states.add(to_state)
88        if from_state in self.transitions:
89            if to_state in self.transitions[from_state]:
90                self.transitions[from_state][to_state].update(input_symbols)
91            else:
```

```python
92                self.transitions[from_state][to_state] = input_symbols
93            else:
94                self.transitions[from_state] = {to_state: input_symbols}
95
96        def add_transition_from_dict(self, translations: Dict[int, Dict[int, set]]):
97            """
98            :param translations: translations[f][t] = d where f is from state,t in to state,
99                                  d is the dict of states where d[state] = set of input symbols
100           :type translations: dict
101           """
102           for from_state, to_states in translations.items():
103               for to_state, input_symbols in to_states.items():
104                   self.add_transition(from_state, to_state, input_symbols)
105
106       def __repr__(self):
107           """
108           display the information of the automata
109           """
110           trans = ""
111           for from_state, to_states in self.transitions.items():
112               for to_state, symbols in to_states.items():
113                   for char in symbols:
114                       trans += f"\t{from_state}->{to_state} on '{char}'\n"
115               trans += '\n'
116
117           return f"states:\t{self.states}\n" \
118               f"start state:\t{self.start_state}\n" \
119               f"final state:\t{self.final_states}\n" \
120               f"transitions:\n{trans}"
121
122       def rename(self, offset: int) -> None:
123           """change the state name to prevent the conflict
124
125           :param offset: offset the number
126           :type offset: int
127           """
128           self.states = set(i+offset for i in self.states)
129           self.start_state += offset
130           self.final_states = set(i+offset for i in self.final_states)
131
132           # change the transition
133           new_transitions = dict()
134           for from_state, to_states in self.transitions.items():
135               new_transitions[from_state+offset] = dict()
136               for to_state in to_states.keys():
137                   new_transitions[from_state+offset][to_state+offset] = \
138                       self.transitions[from_state][to_state]
```

```python
139
140         self.transitions = new_transitions
141
142     def draw(self, save='temp.pdf',seed:int=None) -> None:
143         """
144         draw the graph
145
146         :param save: save the save path (`reference <https://stackoverflow.com/a/20382152>`_)
147         :type save: str
148         :param seed: the node location random seed
149         :type seed: int
150
151         if you haven't installed network2tikz,
152         you need install it by
153
154         .. code-block:: bash
155
156             pip install -U network2tikz
157
158         """
159         from network2tikz import plot
160         nodes = list(self.states)
161         node_colors = [
162             'green!20' if node not in self.final_states else 'blue!20' for node in self.states]
163         node_colors[nodes.index(self.start_state)] = "red!20"
164         edges = []
165         edge_labels = []
166         for from_state, to_states in self.transitions.items():
167             for to_state, symbols in to_states.items():
168                 edges.append((from_state, to_state))
169                 labels = []
170                 for symbol in symbols:
171                     labels.append(symbol)
172                 edge_labels.append("| ".join(labels))
173
174         plot((nodes, edges), save,
175             # layout="spring_layout",
176             seed=seed,
177             canvas=(10,10),
178             node_label_as_id=True,
179             node_color=node_colors,
180             edge_label=edge_labels,
181             edge_math_mode=True, edge_directed=True, edge_curved=0.2,
182             edge_label_position='left')
183
184     @classmethod
185     def empty_construct(cls):
```

```python
        """construct a empty construct of a automata

        :return: the empty automata
        :rtype: Automata
        """
        return cls.basic_construct(set([r'\epsilon']))

    @classmethod
    def basic_construct(cls, symbol: set):
        """construct NFA with a single symbol

        :param symbol: the symbol
        :type symbol: str
        :return: a NFA
        :rtype: Automata
        """
        basic = Automata(symbol)
        basic.set_start_state(1)
        basic.add_final_states(2)
        basic.add_transition(1, 2, set(symbol))
        return basic

    @staticmethod
    def star_operation(nfa):
        """process the star operation

        .. note::

            the nfa is changed after call the method

        :param nfa: the previous NFA
        :type nfa: Automata
        :return: the new NFA after processing star operation
            that means add two string in the begin state and end state
        :rtype: Automata
        """
        for final_state in nfa.final_states:
            nfa.add_transition(nfa.start_state, final_state,
                               set([r"\epsilon"]))
            nfa.add_transition(final_state, nfa.start_state,
                               set([r"\epsilon"]))

        return nfa

    @staticmethod
    def concatenation(basic: Automata, addition: Automata) -> Automata:
        """union two Automata
```

```python
        :param basic: this Automata will be changed after union
        :type basic: Automata
        :param addition: This Automata will be deleted after union
        :type addition: Automata
        :return: [description]
        :rtype: Automata
        """
        # to manage the state name conflict
        offset = max(basic.states)
        addition.rename(offset)

        basic.add_transition_from_dict(addition.transitions)
        for pre_final in basic.final_states:
            basic.add_transition(pre_final, addition.start_state,
                            Automata.empty_string)

        basic.final_states = addition.final_states
        del addition
        return basic

    @staticmethod
    def union(basic: Automata, parallel: Automata) -> Automata:
        """handle the regex s|t by union these NFA

        :param basic: the NFA will change after union
        :type basic: Automata
        :param parallel: the NFA will be deleted after union
        :type parallel: Automata
        :return: The new NFA based on `basic`
        :rtype: Automata
        """
        # rename the two graph
        basic.rename(offset=1)
        offset = max(basic.states)
        parallel.rename(offset)

        # update edges
        basic.add_transition_from_dict(parallel.transitions)

        # update the start
        new_start_state = min(basic.states) - 1
        basic.add_transition(new_start_state,
                        basic.start_state, Automata.empty_string)
        basic.add_transition(new_start_state, parallel.start_state,
                        Automata.empty_string)
        basic.set_start_state(new_start_state)
```

```python
280
281          # handle the final states
282          new_final_state = max(parallel.states)+1
283          pre_finals = basic.final_states.union(parallel.final_states)
284          for pre_final in pre_finals:
285              basic.add_transition(
286                  pre_final, new_final_state, Automata.empty_string)
287          basic.final_states = set([new_final_state])
288
289          del parallel
290          return basic
291
292
293  if __name__ == "__main__":
294      def figure_path(s):
295          return f"../docs/figures/{s}.pdf"
296
297      # basic test
298      test = Automata(set('ab'))
299      test.set_start_state(1)
300      test.add_final_states(2)
301      test.add_final_states(2)
302      test.add_transition(1, 2, set(['a', 'b']))
303      test.add_transition(1, 1, set('b'))
304      print(test.transitions)
305      print(test)
306      test.draw('../docs/figures/test_automata.pdf',seed=2) # 2
307      """ output
308      {1: {2: {'a', 'b'}, 1: {'b'}}}
309      states: {1, 2}
310      start state:  1
311      final state:  {2}
312      transitions:
313              1->2 on 'a'
314              1->2 on 'b'
315              1->1 on 'b'
316      """
317      print(test.transitions)
318      test.rename(3)
319      print(test)
320      """output
321      {1: {2: {'a', 'b'}, 1: {'b'}}}
322      states: {4, 5}
323      start state:  4
324      final state:  {5}
325      transitions:
326              4->5 on 'a'
```

```
327                    4->5 on 'b'
328                    4->4 on 'b'
329            """
330
331        # test basic construct
332        test1 = Automata.basic_construct(set(['a']))
333        test1.draw(save="../docs/figures/basic_a.pdf",seed=1)
334        print(test1)
335        """
336        states: {1, 2}
337        start state:  1
338        final state:  {2}
339        transitions:
340                1->2 on 'a'
341        """
342
343        # test star operation
344        test1 = Automata.star_operation(test1)
345        print(test1)
346        test1.draw('../docs/figures/test_star.pdf',seed=1)
347        r"""output
348            states: {1, 2}
349            start state:  1
350            final state:  {2}
351            transitions:
352                    1->2 on '\epsilon'
353                    1->2 on 'a'
354
355                    2->1 on '\epsilon'
356        """
357
358        # test link operation
359        test1 = Automata.basic_construct(set(['a']))
360        test2 = Automata.basic_construct(set(['b']))
361        test2.draw(figure_path('basic_b'),seed=1)
362        print(Automata.concatenation(test1, test2))
363        test1.draw(save=figure_path('test_concatenation'),seed=2) # 2
364        r"""output
365            states: {1, 2}
366            start state:  1
367            final state:  {2}
368            transitions:
369                    1->2 on 'a'
370                    1->2 on '\epsilon'
371
372                    2->1 on '\epsilon'
373        """
```

```
374
375     # test parallel union
376     test1 = Automata.basic_construct(set(['a']))
377     test2 = Automata.basic_construct(set(['b']))
378     test3 = Automata.union(test1, test2)
379     test3.draw(save=figure_path('test_union'),seed=79744993) # 1111
380     # import random
381     # for i in range(50):
382     #     s = int(random.random() * 100000000)
383     #     test3.draw(f'/tmp/{s}.pdf',seed=s)
384     print(test3)
385     r"""output
386         states: {1, 2, 3, 4, 5, 6}
387         start state:  1
388         final state:  {6}
389         transitions:
390                 2->3 on 'a'
391
392                 4->5 on 'b'
393
394                 1->2 on '\epsilon'
395                 1->4 on '\epsilon'
396
397                 3->6 on '\epsilon'
398
399                 5->6 on '\epsilon'
400     """
```

# 附录 C    正则表达式的解析

```
1     from Automata import Automata
2
3
4     class RegexParser:
5         """
6         store and parse a apttern
7
8         .. code-block:: text
9
10            <regex> ::= <term> '|' <regex>
11                    | <term>
12
13            <term> ::= { <factor> }
14
15            <factor> ::= <base> { '*' }
```

```
16
17          <base> ::= <char>
18                | '\\' <char>
19                | '(' <regex> ')'
20

21
22      :param pattern: the pattern to match the string
23      :type pattern: str
24
25      :ivar self.pattern: the pattern
26      :ivar self.NFA: the NFA machine
27      """
28      # alphabet = set([chr(i) for i in range(65, 91)])\
29      #     .union([chr(i) for i in range(97, 123)])\
30      #     .union([chr(i) for i in range(48, 58)])

31

32
33      def __init__(self, pattern: str):
34          """store and parse a apttern
35
36          :param pattern: the pattern to match the string
37          :type pattern: str
38          """
39          self.pattern = pattern
40          # self.NFA = self.build_NFA()

41
42      # def build_NFA(self):
43      #     """build a NFA from pattern create :class:`Automata.Automata`
44
45      #     :return: the NFA of the current pattern
46
47      #     :rtype: Automata.Automata
48      #     """
49      #     language = set()
50      #     self.buffer = []
51      #     self.automata = []
52      #     previous = r'\epsilon'
53      #     for char in self.pattern:
54      #         if char in self.alphabet:
55      #             pass
56      #             # TODO
57      #     return None

58
59      def peek(self) -> str:
60          """returns the next item of input without consuming it;
61
62          :return: the next character
```

```python
            :rtype: str
            """
            return self.pattern[0]


    def eat(self, item:str) -> None:
        """eat(item) consumes the next item of input, failing if not equal to item.



        :param item: the next item
        :type item: str
        :raises RuntimeError: get the wrong letter.
        """
        if(self.peek() == item):
            self.pattern = self.pattern[1:]
        else:
            raise RuntimeError(f"expect: {item}; got {self.peek()}")

    def next(self) -> str:
        """returns the next item of input and consumes it;

        :return: the next character
        :rtype: str
        """
        c = self.peek()
        self.eat(c)
        return c

    def parse_base_part(self) -> Automata:
        """check the cases encountered

        .. code-block:: text

            <base> ::= <char>
                    |  '\\' <char>
                    |  '(' <regex> ')'

        :return: Automata of this part
        :rtype: Automata
        """
        if self.peek() == '(':
            self.eat('(')
            r = self.parse_regex()
            self.eat(')')
            return r

        elif self.peek() == '\\':
            self.eat('\\')
```

```python
110            esc = self.next()
111            return Automata.basic_construct(esc)
112        else:
113            return Automata.basic_construct(self.next())

114
115    def parse_factor_part(self) -> Automata:
116        base = self.parse_base_part()
117
118        while(self.pattern and self.peek() == '*'):
119            self.eat('*')
120            base = Automata.star_operation(base)
121
122        return base
123
124    def parse_term_part(self) -> Automata:
125        """check that it has not reached the boundary of a term or the end of the input:
126
127        .. code-block:: text
128
129            <term> ::= { <factor> }
130
131        :return: the NFA of this part
132        :rtype: Automata
133        """
134        factor = Automata.empty_construct()
135        while(self.pattern and self.peek() != ')' and self.peek() != '|'):
136            next_factor = self.parse_factor_part()
137            factor = Automata.concatenation(factor, next_factor)
138
139        return factor
140
141    def parse_regex(self) -> Automata:
142        """For regex() method, we know that we must parse at least one term,
143        and whether we parse another
144
145        .. code-block::text
146
147            <regex> ::= <term> '|' <regex>
148                     |  <term>
149
150        :return: the NFA
151        :rtype: Automata
152        """
153        term = self.parse_term_part()
154        if(self.pattern and self.peek() == '|'):
155            self.eat('|')
156            regex = self.parse_regex()
```

```
157                return Automata.union(term, regex)
158            else:
159                return term
160
161
162
163    if __name__ == "__main__":
164        def figure_path(s):
165            return f"../docs/figures/{s}.pdf"
166
167        # test the cases of the only letter
168        test1 = RegexParser("a")
169        print(test1.parse_base_part())
170        """output
171            states: {1, 2}
172            start state:  1
173            final state:  {2}
174            transitions:
175                    1->2 on 'a'
176        """
177
178        # test the escape symbol
179        test2 = RegexParser("\*")
180        print(test2.parse_base_part())
181        """output
182            states: {1, 2}
183            start state:  1
184            final state:  {2}
185            transitions:
186                    1->2 on '*'
187        """
188
189        # test parse factor part
190        test3 = RegexParser('a*')
191        print(test3.parse_factor_part())
192        r"""output
193            states: {1, 2}
194            start state:  1
195            final state:  {2}
196            transitions:
197                    1->2 on '\epsilon'
198                    1->2 on 'a'
199
200                    2->1 on '\epsilon'
201        """
202
203        test4 = RegexParser('ab')
```

```
204    print(test4.parse_term_part())
205    r"""output
206        states: {1, 2, 3, 4, 5, 6}
207        start state:  1
208        final state:  {6}
209        transitions:
210                1->2 on '\epsilon'
211
212                3->4 on 'a'
213
214                2->3 on '\epsilon'
215
216                5->6 on 'b'
217
218                4->5 on '\epsilon'
219    """
220    nfa1 = RegexParser('(a|b)').parse_regex()
221    print(nfa1)
222    # # which is best
223    # import random
224    # import os
225    # import time
226    # l = []
227    # for i in range(100):
228    #     s = int(random.random() * 100000000)
229    #     nfa1.draw(f'/tmp/ab/{s}.pdf',seed=s)
230    #     l.append(s)
231    # time.sleep(10)
232    # for s in l:
233    #     os.system(f"pdftoppm /tmp/ab/{s}.pdf /tmp/ab/{s} -png")
234    nfa1.draw(save=figure_path("a|b"),seed=79870681)
235
236    # test a complex
237    nfa2 = RegexParser('(a|b)*ab').parse_regex()
238    nfa2.draw(save=figure_path('complex'),seed=53138909)
```