

Regular Expressions

Regular Expressions

A notation to specify a set of strings

pattern matching---find one of a specific set of strings

Applications

Test if a string matches some pattern.

- Process natural language.
- Scan for virus signatures.
- Specify a programming language.
- Access information in digital libraries.
- Search genome using PROSITE patterns.
- Filter text (spam, NetNanny, Carnivore, malware).
- Validate data-entry fields (dates, email, URL, credit card).

...

Parse text files.

- Compile a Java program.
- Crawl and index the Web.
- Read in data stored in ad hoc input file format.
- Create Java documentation from Javadoc comments.

REs are amazingly powerful and expressive, but using them in applications can be amazingly complex and error-prone

REs and NFAs

Machine to recognize whether a given string is in a given text

Kleene's theorem

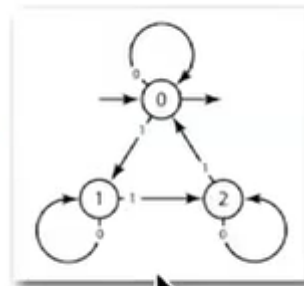
- For any DFA, there exists a RE that describes the same set of strings
- For any RE, there exists a DFA that recognizes the same set of things

RE

$0^* \mid (0^*10^*10^*10^*)^*$

number of 1's is a multiple of 3

DFA



number of 1's is a multiple of 3

NFA

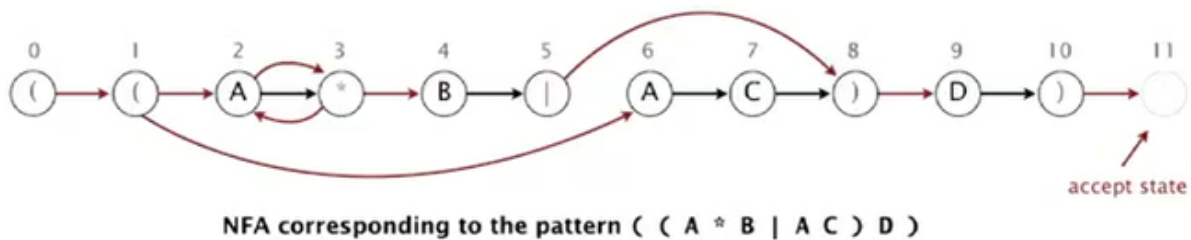
DFA---deterministic finite state automata, easy because exactly one applicable transition

NFA---nondeterministic finite state automata, can be several applicable transitions and need to select the right one

- Nondeterminism
 - one view---machine can guess the proper sequence of state transitions
 - another---sequence is a proof that the machine accepts the text

RE matching NFA

- RE enclosed in parentheses
- One state per RE character (start=0, accept=M)
- Red ϵ -transition (change state, but don't scan text)
- Black match transition (change state and can to next text char)
- Accept if any sequence of transitions ends in accept state



NFA Simulation

Elements

- state names---integers from 0 to M
- match transitions---keep RE in array $re[]$
- ϵ -transitions---store in a digraph G

Simulation

maintain a set of all possible states that NFA could be in after reading in the first i text characters

Complexity

Proposition. Determining whether an N -character text is recognized by the NFA corresponding to an M -character pattern takes time proportional to MN in the worst case.

Pf. For each of the N text characters, we iterate through a set of states of size no more than M and run DFS on the graph of ϵ -transitions.

[The NFA construction we will consider ensures the number of edges $\leq 3M$.]

Digraph reachability

find all vertices reachable from a given source or set of vertices

```
1 public class NFA {
2
3     private Digraph graph;    // digraph of epsilon transitions
4     private String regexp;    // regular expression
5     private final int m;      // number of characters in regular expression
6
7     public NFA(String regexp) {}
8
9     public boolean recognizes(String txt) {
10         DirectedDFS dfs = new DirectedDFS(graph, 0);
11         Bag<Integer> pc = new Bag<Integer>();
12         for (int v = 0; v < graph.V(); v++)
13             if (dfs.marked(v)) pc.add(v);
14
15         // Compute possible NFA states for txt[i+1]
16         for (int i = 0; i < txt.length(); i++) {
17             if (txt.charAt(i) == '*' || txt.charAt(i) == '|' || txt.charAt(i)
18 == '(' || txt.charAt(i) == ')')
19                 throw new IllegalArgumentException("text contains the
20 metacharacter '" + txt.charAt(i) + "'");
21
22             Bag<Integer> match = new Bag<Integer>();
23             for (int v : pc) {
24                 if (v == m) continue;
25                 if ((regexp.charAt(v) == txt.charAt(i)) || regexp.charAt(v) ==
26 '.')
27                     match.add(v+1);
28             }
29             if (match.isEmpty()) continue;
30
31             dfs = new DirectedDFS(graph, match);
32             pc = new Bag<Integer>();
33             for (int v = 0; v < graph.V(); v++)
34                 if (dfs.marked(v)) pc.add(v);
35
36             // optimization if no states reachable
37             if (pc.size() == 0) return false;
38         }
39         return true;
40     }
41 }
```

```

35     }
36
37     // check for accept state
38     for (int v : pc)
39         if (v == m) return true;
40     return false;
41 }
42
43 public static void main(String[] args) {
44     String regexp = "(" + args[0] + ")";
45     String txt = args[1];
46     NFA nfa = new NFA(regexp);
47     Stdout.println(nfa.recognizes(txt));
48 }
49
50 }

```

NFA Construction

Stack

Challenges. Remember left parentheses to implement closure and or; remember | to implement or.

Solution. Maintain a stack.

- (symbol: push (onto stack.
- | symbol: push | onto stack.
-) symbol: pop corresponding (and possibly intervening |; add ϵ -transition edges for closure/or.

Implementation

```

1     public NFA(String regexp) {
2         this.regexp = regexp;
3         m = regexp.length();
4         Stack<Integer> ops = new Stack<Integer>();
5         graph = new Digraph(m+1);
6         for (int i = 0; i < m; i++) {
7             int lp = i;
8             if (regexp.charAt(i) == '(' || regexp.charAt(i) == '|')
9                 ops.push(i);
10            else if (regexp.charAt(i) == ')') {
11                int or = ops.pop();
12
13                // 2-way or operator
14                if (regexp.charAt(or) == '|') {
15                    lp = ops.pop();

```

```

16         graph.addEdge(lp, or+1);
17         graph.addEdge(or, i);
18     }
19     else if (regexp.charAt(or) == '(')
20         lp = or;
21     else assert false;
22 }
23
24 // closure operator (uses 1-character lookahead)
25 if (i < m-1 && regexp.charAt(i+1) == '*') {
26     graph.addEdge(lp, i+1);
27     graph.addEdge(i+1, lp);
28 }
29 if (regexp.charAt(i) == '(' || regexp.charAt(i) == '*' ||
regexp.charAt(i) == ')')
30     graph.addEdge(i, i+1);
31 }
32 if (ops.size() != 0)
33     throw new IllegalArgumentException("Invalid regular expression");
34 }

```

Applications

Completion

To complete the implementation:

- Add wildcard.
- Add multiway or.
- Handle metacharacters.
- Support character classes.
- Add capturing capabilities.
- Extend the closure operator.
- Error checking and recovery.
- Greedy vs. reluctant matching.

Generalized RE Print

Take a RE as a command-line argument and print the lines from standard input having some substring that is matched by the RE

```

1 public class GREP {
2
3     // do not instantiate
4     private GREP() { }
5
6     public static void main(String[] args) {
7         String regexp = "(" + args[0] + ".*)";
8         NFA nfa = new NFA(regexp);

```

```

9         while (StdIn.hasNextLine()) {
10             String line = StdIn.readLine();
11             if (nfa.recognizes(line)) {
12                 StdOut.println(line);
13             }
14         }
15     }
16 }

```

Java string library

java.util.regex.Pattern

java.util.regex.Matcher

```

1  /*****
2  * % java Validate "..oo..oo." bloodroot
3  * true
4  *
5  * % java Validate "..oo..oo." nincompophood
6  * false
7  *
8  * % java Validate "[^aeiou]{6}" rhythm
9  * true
10 *
11 * % java Validate "[^aeiou]{6}" rhythms
12 * false
13 *****/
14
15 public class Validate {
16
17     public static void main(String[] args) {
18         String regexp = args[0];
19         String text = args[1];
20         StdOut.println(text.matches(regexp));
21     }
22 }

```

Harvesting information

print all substrings of input that match a RE

Not-so-regular expressions

pattern matching with back-references is intractable

Back-references.

- `\1` notation matches subexpression that was matched earlier.
- Supported by typical RE implementations.

```
(.+)\1      // beriberi couscous  
1?$|^(11+?)\1+  // 1111 111111 1111111111
```

Some non-regular languages.

- Strings of the form ww for some string w : beriberi.
- Unary strings with a composite number of 1s: 111111.
- Bitstrings with an equal number of 0s and 1s: 01110100.
- Watson-Crick complemented palindromes: atttcggaat.

Summary

- Abstract machines, languages, and nondeterminism
 - Basis of the theory of computation
 - Intensively studied since the 1930s
 - Basis of programming languages
- Compiler---A program that translates a program to machine code
 - **KMP** string \Rightarrow DFA
 - **grep** re \Rightarrow NFA
 - **javac** java language \Rightarrow java byte code

	KMP	grep	Java
pattern	string	RE	program
parser	unnecessary	check if legal	check if legal
compiler output	DFA	NFA	byte code
simulator	DFA simulator	NFA simulator	JVM

- programmer
 - implement substring search via DFA simulation
 - implement RE pattern matching via NFA simulation
- Theoretician
 - RE is a compact description of a set of strings
 - NFA is an abstract machine equivalent in power of RE
 - DFAs and REs have limitations