

Describing a Pattern

PROSITE. Huge database of protein families and domains.

Q. How to describe a protein motif?

Ex. [signature of the C₂H₂-type zinc finger domain]

- c
- Between 2 and 4 amino acids.
- c
- 3 more amino acids.
- One of the following amino acids: LIVMFYWCX.
- 8 more amino acids.
- H
- Between 3 and 5 more amino acids.
- H



CAASCGGPY**ACGGWAGY**HAGWH

Pattern Matching Applications

Test if a string matches some pattern.

- Process natural language.
- Scan for virus signatures.
- Check style of Java variable names.
- Access information in digital libraries.
- Search-and-replace in a word processors.
- Filter text (spam, NetNanny, ads, Carnivore, malware).
- Validate data-entry fields (dates, email, URL, credit card).
- Search for markers in human genome using PROSITE patterns.

MarkovModel.java:16:23:
Name 'NUM_CHAR' must match pattern
'^[a-z][_a-zA-Z0-9]*\$|^[A-Z]\$'.

Parse text files.

- Compile a Java program.
- Crawl and index the Web.
- Read in data stored in TOY input file format.
- Automatically create Java documentation from Javadoc comments.

5

6

Regular Expressions: Basic Operations

Regular expression. Notation to specify a set of strings.

operation	regular expression	matches	does not match
concatenation	aabaab	aabaab	every other string
wildcard	.u.u.u.	cumulus jugulum	succubus tumultuous
union	aa baab	aa baab	every other string
closure	ab*a	aa abbba	ab ababa
parentheses	a(a b)aab	aaaab abaab	every other string
	(ab)*a	a ababababa	aa abbba

Regular Expressions: Examples

Regular expression. Notation is surprisingly expressive.

regular expression	matches	does not match
. *spb. * contains the trigraph spb	raspberry crispbread	subspace subspecies
a* (a*ba*ba*ba*) * multiple of three b's	bbb aaa bbbaababbaa	b bb baabbaa
. *0. . . . fifth to last digit is 0	1000234 98701234	111111111 403982772
gcg(cgg agg)*ctg fragile X syndrome indicator	gcgctg gcgcgctg gcgcgaggctg	gcgcgg cggcgcgctg gcgcgagctg

7

8

Generalized Regular Expressions

Regular expressions are a standard programmer's tool.

- Built in to Java, Perl, Unix, Python,
- Additional operations typically added for convenience.
- Ex: `[a-e]+` is shorthand for `(a|b|c|d|e)(a|b|c|d|e)*`.

operation	regular expression	matches	does not match
one or more	<code>a(bc)+de</code>	abcde abcbcde	ade bcde
character class	<code>[A-Za-z][a-z]*</code>	lowercase Capitalized	camelCase 4illegal
exactly k	<code>[0-9]{5}-[0-9]{4}</code>	08540-1321 19072-5541	111111111 166-54-1111
negation	<code>[^aeiou]{6}</code>	rhythm	decade

Regular Expressions in Java

Validity checking. Is `input` in the set described by the `re`?

```
public class Validate {
    public static void main(String[] args) {
        String re = args[0];
        String input = args[1];
        StdOut.println(input.matches(re));
    }
}
```

powerful string library method

```
% java Validate "C.{2,4}C...[LIVMFYWC].{8}H.{3,5}H" CAASCGGPYACGGAAGYHAGAH
true
% java Validate "[$_A-Za-z][$_A-Za-z0-9]*" ident123
true
% java Validate "[a-z]+@[a-z]+\.(edu|com)" wayne@cs.princeton.edu
true
```

Annotations:
 - `C.{2,4}C...`: C_2H_2 type zinc finger domain
 - `[$_A-Za-z]`: legal Java identifier
 - `[a-z]+@[a-z]+\.(edu|com)`: valid email address (simplified)
 - `" "`: need quotes to "escape" the shell

9

10

String Searching Methods

String Searching Methods

```
public class String (Java's String library)
```

`boolean matches(String re)` *does this string match the given regular expression*

`String replaceAll(String re, String t)` *replace all occurrences of regular expression with the replacement string*

`int indexOf(String t, int from)` *return the index of the first occurrence of the string t after the index from*

`String[] split(String re)` *split the string around matches of the given regular expression*

```
public class String (Java's String library)
```

`boolean matches(String re)` *does this string match the given regular expression*

`String replaceAll(String re, String t)` *replace all occurrences of regular expression with the replacement string*

`int indexOf(String t, int from)` *return the index of the first occurrence of the string t after the index from*

`String[] split(String re)` *split the string around matches of the given regular expression*

```
String s = StdIn.readAll();
s = s.replaceAll("\\s+", " ");
```

replace all sequences of whitespace characters with a single space

```
String s = StdIn.readAll();
String[] words = s.split("\\s+");
```

create array of words in document

regular expression that matches any whitespace character

11

12

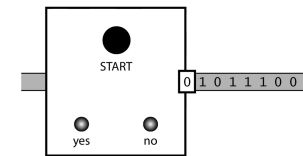
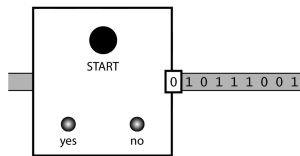
DFAs

Regular expressions are a concise way to describe patterns.

- How would you implement the method `matches()` ?
- Hardware: build a deterministic finite state automaton (DFA).
- Software: simulate a DFA.

DFA: simple machine that solves a pattern match problem.

- Different machine for each pattern.
- Accepts or rejects string specified on input tape.
- Focus on `true` or `false` questions for simplicity.

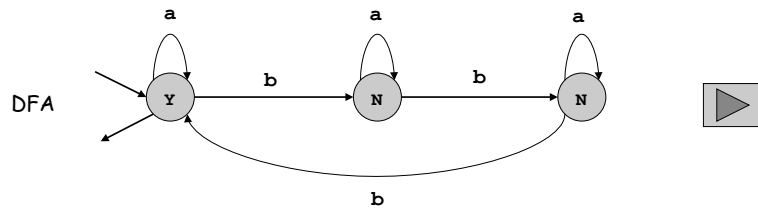


Deterministic Finite State Automaton (DFA)

DFA and RE Duality

Simple machine with N states.

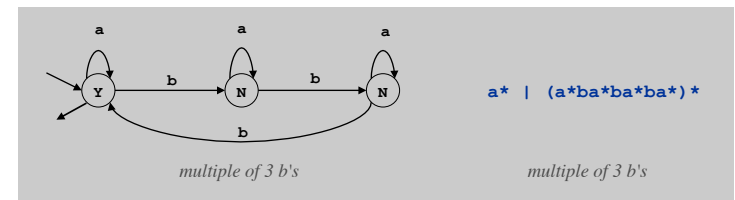
- Begin in **start state** (entering arrow).
- Read first input symbol.
- Move to next state, depending on current state and input symbol.
- Repeat until last input symbol read.
- Accept if last state is an **accept state** (leaving arrow).



RE. Concise way to **describe** a set of strings.

DFA. Machine to **recognize** whether a given string is in a given set.

Duality. 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.



Practical consequence of duality proof: to match RE, (i) build DFA and (ii) simulate DFA on input string.

Implementing a Pattern Matcher

Problem. Given a RE, create program that tests whether given input is in set of strings described.

Step 1. Build the DFA.

- A compiler!
- See COS 226 or COS 320.

Step 2. Simulate it with given input.

```
State state = start;
while (!StdIn.isEmpty()) {
    char c = StdIn.readChar();
    state = state.next(c);
}
StdOut.println(state.accept());
```

17

Application: Harvester

Harvest information from input stream.

- Harvest patterns from DNA.

```
% java Harvester "gcg(cgg|agg)*ctg" chromosomeX.txt
gcgcgggcgggcgggcgggctg
gcgctg
gcgctg
gcgcgggcgggcgggcgggcgggctg
```

- Harvest email addresses from web for spam campaign.

```
% java Harvester "[a-z]+@[a-z]+\.(edu|com)" http://www.princeton.edu/~cos126
rs@cs.princeton.edu
maia@cs.princeton.edu
doug@cs.princeton.edu
wayne@cs.princeton.edu
```

18

Application: Harvester

Harvest information from input stream.

- Use `Pattern` data type to compile regular expression to NFA.
- Use `Matcher` data type to simulate NFA.

equivalent, but more efficient representation of a DFA

```
import java.util.regex.Pattern;
import java.util.regex.Matcher;

public class Harvester {
    public static void main(String[] args) {
        String re = args[0];
        In in = new In(args[1]);
        String input = in.readAll();
        Pattern pattern = Pattern.compile(re);
        Matcher matcher = pattern.matcher(input);

        while (matcher.find()) {
            StdOut.println(matcher.group());
        }
    }
}
```

Annotations in the code block:

- create NFA from RE (pointing to `Pattern.compile(re)`)
- create NFA simulator (pointing to `pattern.matcher(input)`)
- look for next match (pointing to `matcher.find()`)
- the match most recently found (pointing to `matcher.group()`)

19


Application: Parsing a Data File

Ex: parsing an NCBI genome data file.

```
LOCUS AC146846 128142 bp DNA linear HTG 13-NOV-2003
DEFINITION Ornithorhynchus anatinus clone CLM1-393H9,
ACCESSION AC146846
VERSION AC146846.2 GI:38304214
KEYWORDS HTG; HTGS_PHASE2; HTGS_DRAFT.
SOURCE Ornithorhynchus anatinus (platypus)
ORIGIN
1 tgtatttcat ttgacogtgc tgttttttcc cggtttttca gtacgggtgt agggagccac
61 gtgattctgt ttgtttatg ctgcogaata gctgctgat gaatctctgc atagacagct // a comment
121 gccgcagggg gaaatgacca gtttgatgat acaaatgta ggaaagctgt ttctcataa
...
128101 ggaaatgcca cccccacgct aatgtacagc ttcttagat tg
```

Annotations in the code block:

- header info (pointing to the first few lines)
- line numbers (pointing to 1, 61, 121, and 128101)
- spaces (pointing to the space between 1 and tgtatttcat)
- the DNA (pointing to the DNA sequence line)
- comments (pointing to `// a comment`)



20

Ex: parsing an NCBI genome data file.

```
String re = "[ ]*[0-9]+([actg ]*).*";
Pattern pattern = Pattern.compile(re);
In in = new In(filename);
while (!in.isEmpty()) {
    String line = in.readLine();
    Matcher matcher = pattern.matcher(line);
    if (matcher.find()) {
        String s = matcher.group(1).replaceAll(" ", "");
        // do something with s
    }
}
```

```
LOCUS AC146846 128142 bp DNA linear HTG 13-NOV-2003
DEFINITION Ornithorhynchus anatinus clone CLM1-393H9,
ACCESSION AC146846
VERSION AC146846.2 GI:38304214
KEYWORDS HTG; HTGS_PHASE2; HTGS_DRAFT.
SOURCE Ornithorhynchus anatinus (platypus)
ORIGIN
    1  tgtatttcac  ttgaccgtgc  tgttttttcc  cggtttttca  gtaacggtgtt  agggagccac
    61  gtgattctgt  ttgttttatg  ctgccgaata  gctgctcgat  gaatctctgc  atagacagot // a comment
    121 gccgcaggga  gaaatgacca  gtttggatg  acaaaatgta  gaaaagctgt  ttcttcataa
    ...
128101 ggaaatgcga  cccccacgct  aatgtacagc  ttcttagat  tg
```

Summary

Programmer.

- Regular expressions are a powerful pattern matching tool.
- Implement regular expressions with finite state machines.

Theoretician.

- RE is a compact description of a set of strings.
- DFA is an abstract machine that solves RE pattern match problem.

You. Practical application of core CS principles.



<http://xkcd.com/208/>

Fundamental Questions

Q. Are there patterns that **cannot** be described by any RE/DFA?

- A. Yes.
- Bit strings with equal number of 0s and 1s.
 - Decimal strings that represent prime numbers.
 - DNA strings that are Watson-Crick complemented palindromes.

Q. Can we extend RE/DFA to describe richer patterns?

- A. Yes.
- Context free grammar (e.g., Java).
 - **Turing machines.**