CSC236 winter 2020, week 9: Formal languages and regular expressions

Recommended reading: Chapter 7 Vassos course notes

Colin Morris colin@cs.toronto.edu http://www.cs.toronto.edu/~colin/236/W20/

March 9, 2020

Outline

Upcoming dates

Iterative correctness wrap-up

Formal languages

Regular expressions

Upcoming stuff

► A2 due Thurs 3pm



- Extra office hours today 2-4pm.
 Extended afficiency
- Extended office hours Wednesday
- Term test 2, next Monday (March 16)
 - ▶ 12:10-13:00 @ EX320
 - ▶ 13:10-14:00 @ EX200

Term test 2

Covers weeks 4-8 (same material as A2). Potential topics for questions:¹

- Devise a recurrence for the runtime of an algorithm
- Use unwinding to find a closed form for a recurrence
- Use Master Theorem to reason about big- Θ of divide-and-conquer recurrences
- Come up with formal specifications for an algorithm
- ► Use induction to prove correctness of a recursive algorithm
- Identify and prove loop invariants
- Use loop invariants to prove partial correctness
- Prove that an iterative algorithm terminates

¹Not exhaustive.

Term test 2

Prefab cheet sheet will be provided. Will have Master Theorem, and possibly more. e.g.

- Geometric series identities $(\sum_{i=0}^{n} 2^{i} = 2^{n+1} 1)$
- ▶ Big- Θ definition
- Brief reminder of 'recipes' for proofs of recursive correctness, partial correctness, termination, etc.

Will be posted to course website at least a couple days before test. (But not necessarily indicative of what questions will be on the test.)

'Clamping' invariants (Example from week 7 quiz v2)

```
1
    def mult(x, y):
         """Pre: x and y are ints. y is non-negative.
2
3
        Post: return x * y
         .....
4
5
     \mathbf{p} = \mathbf{0}
     while y > 0:
6
7
             p += x
8
             v -= 1
9
        return p
```

Working backwards. I want to say that if the loop exits, y will be 0. What loop invariant will allow this?

inv: 12; 20 Welction , JEO

Clamping invariants: another example

$$m_{j} = G_{j} \times b_{j}$$

$$m_{i} = m_{i} \cap (a_{i,j}, b_{j})$$

$$m_{i} = 100 + G_{i}$$

```
def geq(a, b):
1
    """Pre: a and b are positive integers
2
                                    (nv: 1) G_0 - 6 = 9: -6
    Post: return True iff a >= b
3
    11 11 11
4
    while a > 0 and b > 0:
5
                                         3) a; > 0 1 6 > 0
      a -= 1
6
7
      b -= 1
    return b == 0
8
```

What should be true when the loop exits?

it we exit

 $a_i = 0$ or $b_i = 0$

Anti-pattern: invariants involving loop counter

```
def geq(a, b):
1
     """Pre: a and b are positive integers
2
    Post: return True iff a >= b
3
     .....
4
     while a > 0 and b > 0:
5
       a -= 1
6
7
       h -= 1
     return b == 0
8
```

At the end of each iteration j...

Not wrong, but can be simplified.

Related anti-pattern: counting exact number of iterations

Can work, but generally more work than necessary

```
def f(n):
1
      """Pre: n is a positive integer.
2
3
      11 11 11
4
     a = b = 0
   while n > 0:
5
     if n \% 2 == 1:
6
7
      n -= 1
         a += 1
8
9
       else:
          n = n / / 2
10
          b += 1
11
     return (a, b)
12
```

 $m_j = n_j$

How many times will this loop iterate for a given n?

- ▶ log *n*?
- ▶ $\lceil \log n \rceil$?
- ▶ $\lfloor \log n \rfloor$ + number of 1's in binary representation of *n*?

A: who cores

Invariants that follow directly from code

```
def geq(a, b):
1
     """Pre: a and b are positive integers
2
     Post: return True iff a >= b
3
     .....
4
5
     while a > 0 and b > 0:
       a -= 1
6
     h -= 1
7
     return b == 0
8
```

At the end of each iteration j...

► $a_j = a_{j-1} - 1$ for j > 0

This is (mostly) true, but how can we prove it?

(If your loop invariant just says what the code does, you can probably omit it.)

Example correctness proofs

- Writeup of merge correctness posted with lecture slides
- Sample solutions for tutorial exercises and quizzes
- Vassos notes
- A2Q3 appendix

The rest of the course will be spent studying sets of strings (languages). Relevance to computer science?

- Strings can represent any data type.
 - Ultimately, your computer uses strings of 1's and 0's to represent numbers, lists, trees, cat pictures, etc.
- Sets of strings are a convenient way to formalize the concept of a 'problem' in theoretical computer science
 - P and NP are sets of *languages*
 - ► SAT: the set of strings which represent satisfiable propositional formulas

Heading towards big question of theoretical CS: which problems can be solved algorithmically?

- Specifically, we'll be focusing on the question: which problems can be solved with extremely limited RAM?
- Along the way, we'll need to develop abstract mathematical models for problems, and algorithmic processes

(Connection to computation will become much clearer next week.)

Definitions

- Alphabet: a set of symbols (usually finite), denoted by Σ
 - e.g. $\Sigma = \{0,1\}$, $\Sigma = \{a, b, c, \dots, z\}$
 - \blacktriangleright we'll generally avoid alphabet symbols that could cause confusion, such as $\varepsilon,$ *, parentheses, spaces, etc.

E9 = G

- $\blacktriangleright\ \Sigma^*$ is the set of all finite strings over alphabet Σ
 - e.g. $\{a, b\}^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \ldots\}$
 - recall ε is the **empty string** (like "" in Python)
- $L \subseteq \Sigma^*$ is a **language**
 - ▶ languages may be finite, e.g. $L = \{baa, baabaa\}$
 - ▶ ... or infinite, e.g. $L = \{s \in \{0,1\}^* \mid s \text{ has same number of 0's as 1's}\}$
 - but the strings they contain are always finite
 - NB: $\{\} \neq \{\varepsilon\}$

String operations

Let s, t be strings over some alphabet Σ .

- st is the **concatenation** of s and t
 - occasionally also written $s \circ t$
- ▶ *sⁿ* is repeated concatenation. Defined recursively:

$$(96)^2 = 9696$$

► $s^{j+1} = ss^j$

 $\blacktriangleright s^0 = \varepsilon$

- ▶ |s| is the number of symbols in *s*. Note that $|\varepsilon| = 0$.
 - (len(s) is fine too)

Operations on languages

 $I \cup I'$: union

- $I \cap I'$: intersection
- I I': difference
 - \overline{L} : Complement of L, i.e. $\Sigma^* L$. If L is language of strings over $\{0,1\}$ that start with 0, then \overline{L} is the language of strings that begin with 1 plus the empty string.
 - Sa, 63 20,13 -*LL'*: concatenation, i.e. $\{st \mid s \in L, t \in L'\}$. (What happens when one of these is {} or { ε }?) Z q O, q 1, 60, 61 }

 L^k : concatenation of L with itself k times. $L^0 = \{\varepsilon\}$.



Kleene star A few equivalent definitions

Equ, 66 3 = SE, an, 66,

aaaa, 9966, 6666,}

 L^{\ast} contains the strings formed by concatenating zero or more (not necessarily distinct) strings from L

 $L^* = L^0 \cup L^1 \cup L^2 \cup \ldots$

Define L^* recursively as the smallest set such that:

1. $\varepsilon \in L^*$ 2. if $s \in L$, then $s \in L^*$ 3. if $s \in L$ and $t \in L^*$, then $st \in L^*$ Example: Going from formal description to intuition

Let $L = \{aa, b\}$ Define EVENA = { $s \in \{a, b\}^* \mid s$ has an even number of a's} EVENA $\stackrel{?}{=} L^*$ abq e EVENA L's " strings where q's come in pairs"

From intuition to formal description

) e BIN

Let BIN $\subseteq \{0,1\}^*$ be the language of binary numbers (with no redundant leading zeros).

Give a formal definition of BIN using only:

- One or more finite languages
- > Operations such as union, complement, concatenation, Kleene star, etc.

From intuition to formal description

Let UNIFORM = { $s \in \{a, b\}^* | s$ consists of non-zero repetitions of a single symbol}. Give a formal definition of UNIFORM using only:

- One or more finite languages
- > Operations such as union, complement, concatenation, Kleene star, etc.

えるうちの3* UEを63* $(\xi_{9})^{*} \cup \xi_{6}^{*}) - \xi_{e}$

Regular expressions

Same idea as before, more concise notation

BIN: $0 + (1(0+1)^*)$

UNIFORM:
$$(aa^*) + (bb^*)$$

Not you...

How to validate an email address using a regular expression?

Asked 11 years 4 months ago Active 1 month ago Viewed 1.2m times

-

Over the years I have slowly developed a regular expression that validates MOST email addresses correctly, assuming they don't use an IP address as the server part.

- 3256 I use it in several PHP programs, and it works most of the time. However, from time to time I get contacted by someone that is having trouble with a site that uses it, and I end up having to make some adjustment (most recently I realized that I wasn't allowing 4-character TLDs).
- + What is the best regular expression you have or have seen for validating emails?
- I've seen several solutions that use functions that use several shorter expressions, but I'd rather have one long complex expression in a simple function instead of several short expression in a more complex function.

regex validation email email-validation string-parsing

The fully RFC 822 compliant regex is inefficient and obscure because of its length. Fortunately, RFC 822 was superseded twice and the current specification for email addresses is RFC 5322. RFC 5322 leads to a regex that can be understood if studied for a few minutes and is efficient enough for 2397 actual use -

One RFC 5322 compliant regex can be found at the top of the page at http://emailregex.com/ but uses the IP address pattern that is floating around the internet with a bug that allows 00 for any of the unsigned byte decimal values in a dot-delimited address, which is illegal. The rest of it appears to be consistent with the RFC 5322 grammar and passes several tests using grep -Po, including cases domain names. IP addresses, bad ones, and account names with and without quotes.

Correcting the 00 bug in the IP pattern, we obtain a working and fairly fast regex. (Scrape the rendered version, not the markdown, for actual code.)

(?:[a-z0-9!#\$%&**+/=?^ `{|}~-]+(?:\.[a-z0-9!#\$%&**+/=?^ `{|}~-]+)*|"(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21\x23-\x5b\x5d-\x7fl|\\[\x01-\x09\x0b\x0c\x0e-\x7fl]*")@(?:(?:[a-z0-9](?:[a-z0-9-]*[a-z0-9])? \.)+[a-z0-9](?:[a-z0-9-]*[a-z0-9])?]\[(?:(?:(2(5[0-5])[0-4][0-9])[1[0-9][0-9]][1-9]?[0-9])).)(3)(?:(2(5[0-5])[0-4][0-9])[1-9]?[0-9]]). 5||[0-4][0-9])|1[0-9][0-9]][1-9]?[0-9])|[a-z0-9-]*[a-z0-9]:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7fll\\[/x01-\x09\x0b\x0c\x0e-\x7fl)+)\])

(a-2)

(?)

15

Theory vs. practice

The regular expressions we'll be studying have a close connection to 'regular expressions' in software (e.g. grep, or Python's re library), but there are significant differences.

- Software implementations come packed with *lots* of extra features and syntax (many different dialects)
 - Docs for Python's re library \approx 9,000 words
 - Useful for programmers, bad for theoreticians
 - Some extra features increase expressive power, allowing matching languages which aren't truly 'regular' (more on what this means later)

Our RE syntax will be *very* simple

RE syntax

EU, 13

Recursive definition of \mathcal{RE} , the set of regular expressions over some alphabet Σ :

1.
$$\emptyset, \varepsilon \in \mathcal{RE}$$
 e Z

- 2. every symbol in Σ is in \mathcal{RE}
- 3. if T and S are REs, then so are:
 - 3.1 (T + S) (union) lowest precedence operator
 - 3.2 (TS) (concatenation) middle precedence operator
 - 3.3 T^* (star) highest precedence

The precedence rules allow us to make our REs more readable. e.g. we can write

RE semantics

 $\mathcal{L}(R)$ denotes the language represented by regex R. Base cases:

 $\blacktriangleright \mathcal{L}(\underline{\emptyset}) = \emptyset \quad \{ \}$

(rarely used, but needed for completeness)

- $\blacktriangleright \mathcal{L}(\varepsilon) = \{\varepsilon\}$
- $\mathcal{L}(a) = \{a\}$

 $\widetilde{m{
m {\scriptstyle F}}}$ where a is an arbitrary length-one string from our alphabet Σ

 $-4(01) = \{0,1\}$ = $\{0,1\} = \{0,1\}$

(00)* = {E, 00, 00 00, ...}

Constructor cases. For regular expressions S, T:

- $\blacktriangleright \mathcal{L}(S+T) = \mathcal{L}(S) \cup \mathcal{L}(T) \qquad (STT)$
 - 'take either S or T'

$$\mathcal{L}(ST) = \mathcal{L}(S) \mathcal{L}(T) \prec$$

$$\blacktriangleright \mathcal{L}(T^*) = \mathcal{L}(T)^*$$

'0 or more repetitions of T'

RE identities

Most of these follow from definition. Some require proof. See 7.2.4 in Vassos notes.

- communitativity of union: $R + S \equiv S + R$
- ▶ associativity of union: $(R + S) + T \equiv R + (S + T)$

- associativity of concatenation: $(RS)T \equiv R(ST)$
- left distributivity: $R(S + T) \equiv RS + RT$
- right distributivity: $(S + T)R \equiv SR + TR$
- identity for union: $R + \emptyset \equiv R$
- identity for concatenation: $R\varepsilon \equiv R \equiv \varepsilon R$
- annihilator for concatenation: $\emptyset R \equiv \emptyset \equiv R \emptyset$
- idempotence of Kleene star: $(R^*)^* \equiv R^*$

Examples revisited

$$\underbrace{BIN:}_{0+(1(0+1)^*)} \qquad \underbrace{5 5}_{0} \\ \underbrace{5$$

More examples (600 + 00(+010 + 0(1 + 100 + 10) + 10 + 11)) = L3

$$L\left((0+1)(0+1)(0+1)\right) = L3$$

Devise REs (over $\Sigma = \{0,1\}$) that represent $(0 + 0 + 0) = \{0,1\}$

L → all strings of length 3

 $\sqrt{20.13^3} = L3$

(but not RE)

- L_{3} \rightarrow all strings of length 3 or 4
 - strings that start and end with 0
 - 'sorted' strings (0's appear before 1's)
 - \blacktriangleright strings of the form $0^n 1^n$ (sorted and balanced strings)

(0+1)(0+1) (0+1) + (0+1)(0+1)

 $\sum = \frac{1}{2} \omega_{j} h_{j} \sigma_{j} \sigma_{j} t_{j} ?$

L34 = (0+1)(0+1)(0+1)(e+(0+1))(0+1)(0+1)(0+1)(0+1)(0+1)(0+1)(0+1)

More examples



 $\partial (G+I)^* O + O$

Devise REs (over $\Sigma=\{0,1\})$ that represent

- all strings of length 3
- all strings of length 3 or 4
- strings that start and end with 0
- ► 'sorted' strings (0's appear before 1's) O^{*} (*
- strings of the form $0^n 1^n$ (sorted and balanced strings)

More examples

Devise REs (over $\Sigma = \{0, 1\}$) that represent

- ► all strings of length 3
- ▶ all strings of length 3 or 4
- strings that start and end with 0
- 'sorted' strings (0's appear before 1's)
- strings of the form 0ⁿ1ⁿ (sorted and balanced strings)

More examples

Devise REs (over $\Sigma = \{0, 1\}$) that represent

- ► all strings of length 3
- ▶ all strings of length 3 or 4
- strings that start and end with 0
- 'sorted' strings (0's appear before 1's)
- strings of the form 0ⁿ1ⁿ (sorted and balanced strings)

For every language L, does there exist an RE R such that $\mathcal{L}(R) = L$?

How would we prove that a given language can't be represented by any RE? We'll need another tool.

Proving RE equivalence $\sum_{i=1}^{n} \frac{1}{2} \leq \frac{1}{2} \leq$

$$a = (a + b)^* a(a + b)^* b(a + b)^*$$

 $b = (a + b)^* ab(a + b)^*$

Show² that S and T are equivalent – i.e. $\mathcal{L}(S) = \mathcal{L}(T)$. Proof sketch: A = B $B \in A$

In general, to prove two sets are equal, I need to show mutual inclusion.

Is L(T) a subset of L(S)? Given a string generated by T, I can also generate it from S by setting the middle $(0+1)^*$ to the empty string.

Is L(S) a subset of L(T)?

By inspection, I can see that L(T) is the set of all strings containing 'ab'.

Therefore, it suffices to show that every string generated by S contains 'ab'.

Let s be an arbitrary string generated by S. Then s is of the form $s_1 a s_2 b s_3$, where s_1, s_2 , and s_3 are each arbitrary strings of a's and b's.

I can show that s always contains substring 'ab' based on the following cases for s_2:

if s_2 is empty

S T

if s_2 starts with b

if s_2 ends with a

if s_2 starts with a and ends with b

 $^{^{2}}$ We could **prove** this from first principles, but it would be tedious. We'll generally ask for only an informal proof for problems like this.