

# 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 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:<sup>1</sup>

- ▶ 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- $\Theta$  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

---

<sup>1</sup>Not exhaustive.

## Term test 2

Prefab cheat 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- $\Theta$  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):
2     """Pre: x and y are ints. y is non-negative.
3     Post: return x * y
4     """
5     p = 0
6     while y > 0:
7         p += x
8         y -= 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?

## Clamping invariants: another example

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

What should be true when the loop exits?

## Anti-pattern: invariants involving loop counter

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

At the end of each iteration  $j$ ...

▶  $a_j = a_0 - j$

▶  $b_j = b_0 - j$

Not *wrong*, but can be simplified.



## Related anti-pattern: counting exact number of iterations

Can work, but generally more work than necessary

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

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

- ▶  $\log n$ ?
- ▶  $\lceil \log n \rceil$ ?
- ▶  $\lceil \log n \rceil + \text{number of 1's in binary representation of } n$ ?

## Invariants that follow directly from code

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

At the end of each iteration  $j$ ...

►  $a_j = a_{j-1} - 1$

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

# Formal languages

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

# Formal languages

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  $\Sigma$ 
  - ▶ e.g.  $\Sigma = \{0, 1\}$ ,  $\Sigma = \{a, b, c, \dots, z\}$
  - ▶ we'll generally avoid alphabet symbols that could cause confusion, such as  $\varepsilon$ ,  $*$ , parentheses, spaces, etc.
- ▶  $\Sigma^*$  is the set of all finite strings over alphabet  $\Sigma$ 
  - ▶ e.g.  $\{a, b\}^* = \{\varepsilon, a, b, aa, ab, ba, bb, aaa, \dots\}$
  - ▶ recall  $\varepsilon$  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  $\Sigma$ .

- ▶  $st$  is the **concatenation** of  $s$  and  $t$ 
  - ▶ occasionally also written  $s \circ t$
- ▶  $s^n$  is repeated concatenation. Defined recursively:
  - ▶  $s^0 = \varepsilon$
  - ▶  $s^{j+1} = ss^j$
- ▶  $|s|$  is the number of symbols in  $s$ . Note that  $|\varepsilon| = 0$ .
  - ▶ ( $\text{len}(s)$  is fine too)

# Operations on languages

$L \cup L'$ : union

$L \cap L'$ : intersection

$L - L'$ : difference

$\bar{L}$ : Complement of  $L$ , i.e.  $\Sigma^* - L$ . If  $L$  is language of strings over  $\{0, 1\}$  that start with 0, then  $\bar{L}$  is the language of strings that begin with 1 plus the empty string.

$LL'$ : concatenation, i.e.  $\{st \mid s \in L, t \in L'\}$ .

(What happens when one of these is  $\{\}$  or  $\{\varepsilon\}$ ?)

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



# Kleene star

A few equivalent definitions

$L^*$  contains the strings formed by concatenating zero or more (not necessarily distinct) strings from  $L$

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

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  $\text{EVENA} = \{s \in \{a, b\}^* \mid s \text{ has an even number of } a\text{'s}\}$

$\text{EVENA} \stackrel{?}{=} L^*$

## From intuition to formal description

Let  $\text{BIN} \subseteq \{0, 1\}^*$  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  $\text{UNIFORM} = \{s \in \{a, b\}^* \mid s \text{ consists of non-zero repetitions of a single symbol}\}$ .

Give a formal definition of UNIFORM using only:

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

# Regular expressions

Same idea as before, more concise notation

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

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



## 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

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

1.  $\emptyset, \varepsilon \in \mathcal{RE}$
2. every symbol in  $\Sigma$  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

$$a + bb^*$$

instead of

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



## RE semantics

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

Base cases:

- ▶  $\mathcal{L}(\emptyset) = \emptyset$ 
  - ▶ (rarely used, but needed for completeness)
- ▶  $\mathcal{L}(\varepsilon) = \{\varepsilon\}$
- ▶  $\mathcal{L}(a) = \{a\}$ 
  - ▶ where  $a$  is an arbitrary length-one string from our alphabet  $\Sigma$

Constructor cases. For regular expressions  $S, T$ :

- ▶  $\mathcal{L}(S + T) = \mathcal{L}(S) \cup \mathcal{L}(T)$ 
  - ▶ 'take *either*  $S$  or  $T$ '
- ▶  $\mathcal{L}(ST) = \mathcal{L}(S) \mathcal{L}(T)$
- ▶  $\mathcal{L}(T^*) = \mathcal{L}(T)^*$ 
  - ▶ '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

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

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

## 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  $0^n 1^n$  (sorted and balanced strings)

## Big question

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

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

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

Show<sup>2</sup> that  $S$  and  $T$  are **equivalent** – i.e.  $\mathcal{L}(S) = \mathcal{L}(T)$ .

---

<sup>2</sup>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.