

Theory of Computation 1

Sets and Regular Expressions

Frank Stephan

Department of Computer Science

Department of Mathematics

National University of Singapore

fstephan@comp.nus.edu.sg

Languages

Examples of Languages (Sets)

- (a) Set of all natural numbers in binary notation: either **0** or **1** followed by arbitrarily many digits from **0, 1**.
- (b) Set of all possible computer programs in syntax of programming language C: Tools can translate a formal description of C syntax into a syntax checker.
- (c) Set of all C programs which pass a compiler without error messages: Compilers check more than just syntactical correctness.
- (d) Set of all C programs which do not have bugs: No computer program can solve this task completely.
- (e) Set of all texts of books written in English and published between 1066 and 1492: Exhaustive list describes this set.

Theory of Computation

How does one describe above sets? How does one modify descriptions? Do descriptions allow membership checks?

Languages

Language = Set of Strings over an Alphabet.

Alphabet Σ , for example $\Sigma = \{0, 1, 2\}$. Always finite.

Finite languages

$L_1 = \emptyset$, no elements.

$L_2 = \{\varepsilon\}$, set consisting of empty string.

$L_3 = \{00, 01, 02, 10, 11, 12, 20, 21, 22\}$, all elements of length 2.

$L_4 = \{\varepsilon, 0, 00, 000, 0000\}$, all strings of 0s up to length 4.

$L_5 = \{01, 001, 02, 002\}$, all strings consisting of one or two 0s followed by a 1 or 2.

Operations with Languages

Union:

$$L \cup H = \{u : u \in L \text{ or } u \in H\};$$

$$\{00, 01, 02\} \cup \{01, 11, 21\} = \{00, 01, 02, 11, 21\};$$

$$\{0, 00, 000\} \cup \{00, 000, 0000\} = \{0, 00, 000, 0000\}.$$

Intersection:

$$L \cap H = \{u : u \in L \text{ and } u \in H\};$$

$$\{00, 01, 02\} \cap \{01, 11, 21\} = \{01\};$$

$$\{0, 00, 000\} \cap \{00, 000, 0000\} = \{00, 000\}.$$

Set Difference:

$$L - H = \{u : u \in L \text{ and } u \notin H\};$$

$$\{00, 01, 02\} - \{01, 11, 21\} = \{00, 02\}.$$

Concatenation:

$$000 \cdot 1122 = 0001122;$$

$$L \cdot H = \{v \cdot w : v \in L \text{ and } w \in H\};$$

$$\{0, 00\} \cdot \{1, 2\} = \{01, 001, 02, 002\}.$$

Kleene Star and Plus

Definition

$$\mathbf{L}^* = \{\varepsilon\} \cup \mathbf{L} \cup \mathbf{L} \cdot \mathbf{L} \cup \mathbf{L} \cdot \mathbf{L} \cdot \mathbf{L} \cup \dots$$

$$= \{\mathbf{w}_1 \cdot \mathbf{w}_2 \cdot \dots \cdot \mathbf{w}_n : n \geq 0 \text{ and } \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n \in \mathbf{L}\};$$

$$\mathbf{L}^+ = \mathbf{L} \cup \mathbf{L} \cdot \mathbf{L} \cup \mathbf{L} \cdot \mathbf{L} \cdot \mathbf{L} \cup \dots$$

$$= \{\mathbf{w}_1 \cdot \mathbf{w}_2 \cdot \dots \cdot \mathbf{w}_n : n > 0 \text{ and } \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n \in \mathbf{L}\}.$$

Examples

$$\emptyset^* = \{\varepsilon\}.$$

Σ^* is the set of all words over Σ .

$$\{0\}^* = \{\varepsilon, 0, 00, 000, 0000, \dots\}.$$

$\{00, 01, 10, 11\}^*$ are all binary words of even length.

$$\varepsilon \in \mathbf{L}^+ \text{ iff } \varepsilon \in \mathbf{L}.$$

Notation

Often \mathbf{w}^* in place of $\{\mathbf{w}\}^*$;

Often $\mathbf{w} \cdot \mathbf{L}$ in place of $\{\mathbf{w}\} \cdot \mathbf{L}$.

Regular Languages

Regular expressions are either finite sets listed by their elements or obtained from other regular expressions by forming the Kleene star, Kleene plus, union, intersection, set-difference or concatenation.

A language is regular iff it can be described by a regular expression.

Regular sets have many different regular expressions.

For example, $\{0, 00\} \cdot \{1, 2\}$ and $\{01, 001, 02, 002\}$ describe the same set. Also 0^* and $(00)^* \cup 0 \cdot (00)^*$ describe the same set.

Intersections and set difference are traditionally not used in regular expressions, as every regular set has an expression only using union, concatenation and Kleene star.

The complement of a language L is $\Sigma^* - L$.

Quiz

Which three sets are described by two of the following regular expressions?

1. $\{00, 000\}^+$;
2. $\{000, 0000\}^+$;
3. $00 \cdot 0^*$;
4. $000 \cdot 0^*$;
5. $\{000, 0000\} \cup (000000 \cdot 0^*)$;
6. $\{00, 01, 02, 10, 11, 12\}$;
7. $0^*1^*2^*$;
8. $(0^*1^*2^*)^*$;
9. $(\{0, 1\} \cdot \{0, 1, 2\}^*) \cap (\{0, 1, 2\} \cdot \{0, 1, 2\})$.

Exercises 1.6, 1.7 and 1.8

Exercise 1.6

Assume A has 3 and B has 2 elements. How many elements do the following sets have at least and at most; it depends on the actual choice which of the bounds is realised: $A \cup B$, $A \cap B$, $A \cdot B$, $A - B$, $A^* \cap B^*$.

Exercise 1.7

Let A, B be finite sets and $|A|$ be the number of elements of A . Is the following formula correct:

$$|A \cup B| + |A \cap B| = |A| + |B|?$$

Prove your answer.

Exercise 1.8

Make a regular expression without intersection and set difference for $0^*1^*0^*1^* \cap (11)^*(00)^*(11)^*(00)^*$.

Tutorial

Register for a tutorial group

During the semester, you can make up to **FOURTEEN** marks (out of **100** marks), 2 marks per homework.

(a) Reserve question in Forum by putting an entry with Tutorial number and Homework number into the headline, for example T2-1.26 for reserving Homework 1.26 for presentation in Tutorial T2. Use "T0" if you cannot register for tutorial. No two students of a tutorial group should select the same homework, check against the Forum title lines before reserving. Reserve at most one homework per week.

(b) Write up exercise in Forum and present the tutorial in the week where tutorials are presented. You can share your write-up over Zoom. Students of different tutorial groups can solve a homework together, but each of them puts a write-up and makes a presentation.

Tests

The **Midterm Test** occurs in 2021 in Week 8 of the lecture, at the usual lecture time. This is Friday 8 October 2021 and students should log into the Zoom at 10:00 hrs for identity checking, the test is from 10:30 to 11:30 hrs. The test type is “Zoom proctoring with Luminus files”.

The final examination counts 60 marks. It is 27 November 2021 at 9:00 hrs and the duration is 2 hours. Please consult NUS webpages for more information and doublecheck the information there.

Summary: Tutorial 14 marks, Midterm 26 marks, Final Exam 60 marks.

<https://nusmods.com/modules/CS3231/theory-of-computation>

Theorem of Lyndon and Schützenberger

Theorem

If two words \mathbf{v}, \mathbf{w} satisfy $\mathbf{vw} = \mathbf{wv}$ then $\exists \mathbf{u} [\mathbf{v}, \mathbf{w} \in \mathbf{u}^*]$;
If all words $\mathbf{v}, \mathbf{w} \in \mathbf{L}$ satisfy $\mathbf{vw} = \mathbf{wv}$ then $\exists \mathbf{u} [\mathbf{L} \subseteq \mathbf{u}^*]$.

Proof

Case $\mathbf{v} = \varepsilon$ or $\mathbf{w} = \varepsilon$: $\mathbf{u} = \mathbf{vw}$.

Case $\mathbf{v} \neq \varepsilon$ and $\mathbf{w} \neq \varepsilon$ and $|\mathbf{v}| = |\mathbf{w}|$: $\mathbf{v} = \mathbf{w}$.

Case $\mathbf{v} \neq \varepsilon$ and $\mathbf{w} \neq \varepsilon$ and $|\mathbf{v}| < |\mathbf{w}|$: Let \mathbf{k} be greatest common divisor of $|\mathbf{v}|, |\mathbf{w}|$ and \mathbf{u} be the first \mathbf{k} symbols of \mathbf{v} .
There are \mathbf{i}, \mathbf{j} with $\mathbf{v} = \mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_i$ and $\mathbf{w} = \mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_j$ for words $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_j$ of length \mathbf{k} .

Note that $\mathbf{v}^j \mathbf{w}^i = \mathbf{w}^i \mathbf{v}^j$ and $|\mathbf{v}^j| = |\mathbf{w}^i|$, as both have the length \mathbf{ijk} . Thus $\mathbf{v}^j = \mathbf{w}^i$.

For each \mathbf{u}_h there is a position of \mathbf{u}_h in \mathbf{w}^i where \mathbf{u}_1 is at the same position in \mathbf{v}^j . Thus $\mathbf{u}_h = \mathbf{u}_1$.

So $\mathbf{v}, \mathbf{w} \in \mathbf{u}^*$ for $\mathbf{u} = \mathbf{u}_1$.

Example, Second Part

Example: Let $v = abcd$ and $w = abcdef$ and $vw = wv$.
Now $k = 2$ (greatest common divisor of 4, 6).

$$v^3 = ab\ cd\ ab\ cd\ ab\ cd;$$

$$w^2 = ab\ cd\ ef\ ab\ cd\ ef.$$

Now $ab = ab$ at 0, 1, $ab = ef$ at 4, 5, $ab = cd$ at 8, 9. So
 $ab = cd = ef$ and $v, w \in (ab)^*$.

Second Part: Let v be shortest nonempty word of L and u
be shortest word with $v \in u^*$.

Let $w \in L$ be arbitrary.

There is \tilde{u} with $v, w \in \tilde{u}^*$.

Now $\tilde{u}^i = u^j = v$ for some i, j .

Thus $\tilde{u}, u \in \hat{u}^*$ as in Part 1 for some \hat{u} where $|\hat{u}|$ is greatest
common divisor of $|u|$ and $|\tilde{u}|$.

By choice of u , $\hat{u} = u$ and $w \in u^*$.

So $L \subseteq u^*$.

Structural Induction

Theorem

Let P be a property of sets such that the following holds:

- Every finite set (including \emptyset) satisfies P ;
- If L, H satisfy P so does $L \cup H$;
- If L, H satisfy P so does $L \cdot H$;
- If L satisfies P so does L^* .

Then all regular sets satisfy P .

Proof.

Let $L(\sigma)$ be the set generated by the regular expression σ . Here finite sets and the operations union, concatenation and Kleene star are permitted for regular expressions.

Now it is shown that there is no shortest regular expression σ such that $L(\sigma)$ does not satisfy P .

Shortest Expression

Assume that σ is a shortest expression not satisfying P ; if there are several shortest ones, σ is just any of these.

- If σ is a list of a finite set then P satisfies $L(\sigma)$;
- If $\sigma = (\rho \cup \tau)$ then ρ, τ are shorter than σ and $L(\rho), L(\tau)$ satisfy P and so does $L(\sigma) = L(\rho) \cup L(\tau)$;
- If $\sigma = (\rho \cdot \tau)$ then ρ, τ are shorter than σ and $L(\rho), L(\tau)$ satisfy P and so does $L(\sigma) = L(\rho) \cdot L(\tau)$;
- If $\sigma = \tau^*$ then τ is shorter than σ and $L(\tau)$ satisfies P and so does $L(\sigma) = (L(\tau))^*$.

So there is no case in which $L(\sigma)$ would not satisfy P , thus this σ does not exist and there is no regular expression σ for which $L(\sigma)$ does not satisfy P . All regular languages satisfy P .

Strengthening the Theorem

Theorem

Let P be a property of sets such that the following holds:

- The empty set \emptyset and the set $\{\epsilon\}$ satisfy P ;
- For every letter a , the set $\{a\}$ satisfies P ;
- If L, H satisfy P so does $L \cup H$;
- If L, H satisfy P so does $L \cdot H$;
- If L satisfies P so does L^* .

Then all regular sets satisfy P .

This strengthening is just based on the fact that every finite set of strings can be formed using concatenation and union from the sets containing a single letter word, the set containing the empty word and the empty set.

Polynomial and Exponential Growth

Definition

A language L has **polynomial growth** iff there is a polynomial p such that for all n there are in L at most $p(n)$ many words shorter than n .

A language L has **exponential growth** iff there are constants h, k such that L contains for all n at least 2^n words which are shorter than $h \cdot n + k$.

Theorem

Every regular set has either polynomial or exponential growth.

This will be proven by structural induction.

Examples

Every finite set has polynomial growth, as one plus the number of elements is a polynomial which is an upper bound as required.

The set 0^*1^* has polynomial growth as there are $n(n+1)/2$ many words shorter than n in this set.

The set $0^* \cup 1^*$ has polynomial growth as there are at most $2n$ many words shorter than n in this set.

The set $\{00, 11\}^* \cdot \{222222\}$ has exponential growth as it has for all n at least 2^n words shorter than $7 + 2n$.

The set $\{0000, 1111\}^*$ has exponential growth as it has for all n at least 2^n words shorter than $1 + 4n$.

The set $\{00, 11\}^* \cdot \emptyset$ is empty and has polynomial growth.

Rules for Growth

Finite sets have polynomial growth.

If L and H have polynomial growth then so do $L \cup H$ and $L \cdot H$.

If L or H have exponential growth then so does $L \cup H$.

The sets $L \cdot \emptyset$ and $\emptyset \cdot L$ have polynomial growth.

If L and H are not empty and at least one of them has exponential growth so does $L \cdot H$.

If L contains v, w with $vw \neq wv$ then L^* has exponential growth else $L^* \subseteq u^*$ for some u and has polynomial growth.

Let $P(L)$ say that the language L has either polynomial or exponential growth. Then the rules imply that all finite sets satisfy P and that, whenever L, H satisfy P so do $L \cup H$, $L \cdot H$ and L^* .

Thus all regular sets satisfy P by structural induction.

Quiz

1. Does $L \cap H$ have exponential growth whenever L and H have exponential growth?
2. Does $\{0101, 010101\}^*$ have exponential growth?
3. Does $\{000, 001, 011, 111\}^* \cdot \{0000, 1111\}$ have exponential growth?
4. Does the (non-regular) set $A = \{w \in \{0, 1\}^* : w \text{ has at most } \log(|w|) \text{ many } 1\text{s}\}$ have polynomial growth?
5. Does the set $A = \{w \in \{0, 1\}^* : w \text{ has at most } \log(|w|) \text{ many } 1\text{s}\}$ from item 4 have exponential growth?
6. Is there a polynomial p such that every regular set has either exponential growth or at most $p(n)$ elements shorter than n for every n ?

Rules for Regular Expressions

- (a) $L \cup L = L$, $L \cap L = L$, $(L^*)^* = L^*$, $(L^+)^+ = L^+$;
- (b) $(L \cup H)^* = (L^* \cdot H^*)^*$ and if $\varepsilon \in L \cap H$ then $(L \cup H)^* = (L \cdot H)^*$;
- (c) $(L \cup \{\varepsilon\})^* = L^*$, $\emptyset^* = \{\varepsilon\}$ and $\{\varepsilon\}^* = \{\varepsilon\}$;
- (d) $L^+ = L \cdot L^* = L^* \cdot L$ and $L^* = L^+ \cup \{\varepsilon\}$;
- (e) $(L \cup H) \cdot K = (L \cdot K) \cup (H \cdot K)$ and $K \cdot (L \cup H) = (K \cdot L) \cup (K \cdot H)$;
- (f) $(L \cup H) \cap K = (L \cap K) \cup (H \cap K)$ and $(L \cap H) \cup K = (L \cup K) \cap (H \cup K)$;
- (g) $(L \cup H) - K = (L - K) \cup (H - K)$ and $(L \cap H) - K = (L - K) \cap (H - K)$.

Inequality Rules

- (a) $\mathbf{L} \cdot \mathbf{L}$ can be different from \mathbf{L} : $\{0\} \cdot \{0\} = \{00\}$;
- (b) $(\mathbf{L} \cap \mathbf{H})^* \subseteq \mathbf{L}^* \cap \mathbf{H}^*$;
Properness: $\mathbf{L} = \{00\}$, $\mathbf{H} = \{000\}$, $(\mathbf{L} \cap \mathbf{H})^* = \{\varepsilon\}$,
 $\mathbf{L}^* \cap \mathbf{H}^* = \{000000\}^*$;
- (c) If $\{\varepsilon\} \cup (\mathbf{L} \cdot \mathbf{H}) = \mathbf{H}$ then $\mathbf{L}^* \subseteq \mathbf{H}$;
Properness: $\mathbf{L} = \{\varepsilon\}$, $\mathbf{H} = \{0\}^*$;
- (d) If $\mathbf{L} \cup (\mathbf{L} \cdot \mathbf{H}) = \mathbf{H}$ then $\mathbf{L}^+ \subseteq \mathbf{H}$;
Properness: $\mathbf{L} = \{\varepsilon\}$, $\mathbf{H} = \{0\}^*$;
- (e) $(\mathbf{L} \cap \mathbf{H}) \cdot \mathbf{K} \subseteq (\mathbf{L} \cdot \mathbf{K}) \cap (\mathbf{H} \cdot \mathbf{K})$;
Properness: $(\{0\} \cap \{00\}) \cdot \{0, 00\} = \emptyset \subset \{000\} =$
 $(\{0\} \cdot \{0, 00\}) \cap (\{00\} \cdot \{0, 00\})$;
- (f) $\mathbf{K} \cdot (\mathbf{L} \cap \mathbf{H}) \subseteq (\mathbf{K} \cdot \mathbf{L}) \cap (\mathbf{K} \cdot \mathbf{H})$;
Properness: $\{0, 00\} \cdot (\{0\} \cap \{00\}) = \emptyset \subset \{000\} =$
 $(\{0, 00\} \cdot \{0\}) \cap (\{0, 00\} \cdot \{00\})$.

Characterising Kleene Star

Corollary 1.17. For any set L , the following statements characterise L^* and L^+ :

- (a) L^* is the smallest set H such that $\{\varepsilon\} \cup (L \cdot H) = H$;
- (b) L^* is the smallest set H such that $\{\varepsilon\} \cup (L \cdot H) \subseteq H$;
- (c) L^+ is the smallest set H such that $L \cup (L \cdot H) = H$;
- (d) L^+ is the smallest set H such that $L \cup (L \cdot H) \subseteq H$.

In the above, one could also use $H \cdot L$ in place of $L \cdot H$.

Exercise 1.18

Which three of the following sets are not equal to any of the other sets:

- (a) $\{01, 10, 11\}^*$;
- (b) $((\{0, 1\} \cdot \{0, 1\}) - \{00\})^*$;
- (c) $(\{01, 10\} \cdot \{01, 10, 11\} \cup \{01, 10, 11\} \cdot \{01, 10\})^*$;
- (d) $(\{01, 10, 11\} \cdot \{01, 10, 11\})^* \cup \{01, 10, 11\} \cdot (\{01, 10, 11\} \cdot \{01, 10, 11\})^*$;
- (e) $\{0, 1\}^* - \{0, 1\} \cdot \{00, 11\}^*$;
- (f) $((\{01\}^* \cup \{10\})^* \cup \{11\})^*$;
- (g) $(\{\varepsilon\} \cup (\{0\} \cdot \{0, 1\}^* \cap \{1\} \cdot \{0, 1\}^*))^*$.

Explain your answer.

Exercise 1.19

Make a regular expression which contains all those decimal natural numbers which start with **3** or **8** and have an even number of digits and end with **5** or **7**.

Make a further regular expression which contains all odd ternary numbers without leading **0**s; here a ternary number is a number using the digits **0, 1, 2** with **10** being three, **11** being four and **1212** being fifty. The set described should contain the ternary numbers

1, 10, 12, 21, 100, 102, 111, 120, 122, 201, ... which are the numbers **1, 3, 5, 7, 9, 11, 13, 15, 17, 19, ...** in decimal.

Exercise 1.20

Let \mathcal{S} be the smallest class of languages such that

- every language of the form u^* for a nonempty word u is in \mathcal{S} ;
- the union of two languages in \mathcal{S} is again in \mathcal{S} ;
- the concatenation of two languages in \mathcal{S} is again in \mathcal{S} .

Prove by structural induction the following properties of \mathcal{S} :

- (a) Every language in \mathcal{S} is infinite;
- (b) Every language in \mathcal{S} has polynomial growth.

Lay out all inductive steps explicitly without only citing results in this lecture.

Exercise 1.21

Let L satisfy the following statement: For all $u, v, w \in L$, either $uv = vu$ or $uw = wu$ or $vw = wv$. Which of the following statements are true for all such L :

- All $x, y \in L$ satisfy $xy = yx$;
- All sufficiently long $x, y \in L$ satisfy $xy = yx$;
- The language L has polynomial growth.

Give an answer to these questions and prove them.

Exercises 1.22-1.23

In the following, each digit, the symbol ε , the symbol \cdot , the symbol \cup , the comma and each set bracket and each normal bracket have length **1** and the length of the expression is the number of all the symbols in it (counting repetitions).

Exercise 1.22

Let **L** consist of all words which contain each of the letters **0, 1, 2, 3** exactly once. Make a regular expression generating **L** which has at most length **100**.

Exercise 1.23

Make a regular expression for the set $\{\mathbf{w} \in \{0\}^* : |\mathbf{w}| \leq \mathbf{9}\}$ which has at most length 26.

Exercises 1.24-1.26

Let **V** be the set of vowels, **W** be the set of consonants and **S** be the set of punctuation marks and **T** be the set of spacings (blanks and new lines and so on).

Exercise 1.24. Make a regular expression (using above sets) of all words which contain at least two vowels and before, after and between vowels is exactly one consonant.

Exercise 1.25. Make a regular expression of all sentences where each sentence consists of words containing one vowel and arbitrarily many consonants and between two words are spacings and after the last word is a punctuation mark.

Exercise 1.26. Make a regular expressions generating texts of sentences separated by spacings where sentences are as above with the only difference that words can have one or two vowels and up to four consonants.

Exercises 1.27-1.30

Let \mathbf{L}, \mathbf{H} be infinite subsets of $\{1, 2, 3, 4, 5, 6, 7, 8, 9\} \cdot \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}^*$ which are viewed as natural numbers having their usual value. Decide whether the following statements are true and either provide a proof or a regular expression as a counterexample.

Exercise 1.27. If \mathbf{L} does not contain numbers of the form $x, x + 1$, the same is true for $\mathbf{L} \cdot \mathbf{L}$.

Exercise 1.28. There is a number $p \geq 2$ and some \mathbf{L} such that both \mathbf{L} and \mathbf{L}^+ consist only of powers of p .

Exercise 1.29. Let \mathbf{L}^q denote the concatenation of q copies of \mathbf{L} . Find p, q such that the following property $(*)$ is true:
 $(*)$: Every infinite set \mathbf{L} has an infinite subset \mathbf{H} such that \mathbf{H}^q consists only of numbers divisible by p .

Exercise 1.30. Prove that $(*)$ is false when $p = 74$.