

CS_275 Automata and Formal Language Theory

Course Notes

Part II: The Recognition Problem (II)

Chapter II.2.: Basics of Regular Languages and Expressions

Anton Setzer

(Based on a book draft by J. V. Tucker and K. Stephenson)

Dept. of Computer Science, Swansea University

[http://www.cs.swan.ac.uk/~csetzer/lectures/
automataFormalLanguage/current/index.html](http://www.cs.swan.ac.uk/~csetzer/lectures/automataFormalLanguage/current/index.html)

March 17, 2017

II.2.1. Regular Languages (12.2)

II.2.2. Regular Expressions (13.8)

II.2.1. Regular Languages (12.2)

II.2.2. Regular Expressions (13.8)

Finite Languages are Regular

grammar	$G^{ab,aabb,aaabbb}$
terminals	a, b
nonterminals	S
start symbol	S
productions	$S \rightarrow ab$ $S \rightarrow aabb$ $S \rightarrow aaabbb$

This grammar is not regular, since there can only be one terminal in the right hand string. But we can amend this:

Finite Languages are Regular

grammar	$G^{ab,aabb,aaabbb}$
terminals	a, b
nonterminals	$S, S_1, S_2, S_3, S_4, S_5, S_6, S_7, S_8, S_9$
start symbol	S
productions	$S \rightarrow aS_1, S_1 \rightarrow b$ $S \rightarrow aS_2, S_2 \rightarrow aS_3, S_3 \rightarrow bS_4, S_4 \rightarrow b$ $S \rightarrow aS_5, S_5 \rightarrow aS_6, S_6 \rightarrow aS_7, S_7 \rightarrow bS_8,$ $S_8 \rightarrow bS_9, S_9 \rightarrow b$

Finite Languages are Regular

If one wishes, the above grammar can of be optimised as follows:

grammar	$G^{ab, aabb, aaabbb}$
terminals	a, b
nonterminals	$S, S_1, S_3, S_4, S_7, S_8, S_9$
start symbol	S
productions	$S \rightarrow aS_1, S_1 \rightarrow b$ $S_1 \rightarrow aS_3, S_3 \rightarrow bS_4, S_4 \rightarrow b$ $S_3 \rightarrow aS_7, S_7 \rightarrow bS_8,$ $S_8 \rightarrow bS_9, S_9 \rightarrow b$

Lemma II.2.1.1.

Lemma (II.2.1.1.)

All finite languages are regular, and a regular grammar for them can be computed.

Proof: Extend the example above.

A Left-Linear Grammar for $a^m b^n$

The following left-linear grammar generates $\{a^m b^n \mid m, n \geq 1\}$.

grammar	$G^{\text{left-linear}, a^m b^n}$
terminals	a, b
nonterminals	S, T
start symbol	S
productions	$S \rightarrow Sb$ $S \rightarrow Tb$ $T \rightarrow Ta$ $T \rightarrow a$

A Right-Linear Grammar for $a^m b^n$

The following right-linear grammar generates $\{a^m b^n \mid m, n \geq 1\}$:

grammar	$G_{\text{right-linear}, a^m b^n}$
terminals	a, b
nonterminals	S, T
start symbol	S
productions	$S \rightarrow aS$ $S \rightarrow aT$ $T \rightarrow bT$ $T \rightarrow b$

Right-Linear Grammar for Numbers

Here is a right-linear grammars for numbers without leading zeros. We use “|” as for BNF.

grammar	G^{Number}
terminals	$0, 1, \dots, 9$
nonterminals	$Number, Digits$
start symbol	$Number$
productions	$Number \longrightarrow 0$ $Number \longrightarrow 1 Digits \mid 2 Digits \mid \dots \mid 9 Digits$ $Digits \longrightarrow 0 Digits \mid 1 Digits \mid \dots \mid 9 Digits$ $Digits \longrightarrow \epsilon$

Right-Linear Grammar for Numbers

Why didn't we use the following as in BNF?

grammar	G^{Number}
terminals	$0, 1, \dots, 9$
nonterminals	$Number, Digit, NonZeroDigit, Digits$
start symbol	$Number$
productions	$Number \longrightarrow Digit \mid NonZeroDigit \ Digits$ $Digits \longrightarrow Digit \mid Digit \ Digits$ $Digit \longrightarrow 0 \mid NonZeroDigit$ $NonZeroDigit \longrightarrow 1 \mid 2 \mid \dots \mid 9$

Answer:

Right-Linear Grammar for Post Codes

The next grammar generates the postcodes of the form SA1 8PP or in general LLd dLL for digits d and capital letters L without any leading zeros. We use the notation $_$ as in BNF. We write $_$ for blank

Right-Linear Grammar for Post Codes

grammar	$G^{Postcode}$
terminals	0, 1, ..., 9, A, B, ..., Z, \sqcup
nonterminals	<i>postcode</i> , <i>letter2</i> , <i>digit1</i> , <i>blank1</i> , <i>digit2</i> , <i>letter3</i> , <i>letter4</i>
start symbol	<i>postcode</i>
productions	$ \begin{aligned} & \textit{postcode} \longrightarrow A \textit{letter2} \mid B \textit{letter2} \mid \dots \mid Z \textit{letter2} \\ & \textit{letter2} \longrightarrow A \textit{digit1} \mid B \textit{digit1} \mid \dots \mid Z \textit{digit1} \\ & \textit{digit1} \longrightarrow 0 \textit{blank1} \mid 1 \textit{blank1} \mid \dots \mid 9 \textit{blank1} \\ & \textit{blank1} \longrightarrow \sqcup \textit{digit2} \\ & \textit{digit2} \longrightarrow 0 \textit{letter3} \mid 1 \textit{letter3} \mid \dots \mid 9 \textit{letter3} \\ & \textit{letter3} \longrightarrow A \textit{letter4} \mid B \textit{letter4} \mid \dots \mid Z \textit{letter4} \\ & \textit{letter4} \longrightarrow A \mid B \mid \dots \mid Z \end{aligned} $

Example Derivation

Here is a derivation of $SA2_L_8PP \in L(G^{Postcode})$:

postcode \Rightarrow *S letter2*
 \Rightarrow *SA digit1*
 \Rightarrow *SA1 blank1*
 \Rightarrow *SA1_L_ digit2*
 \Rightarrow *SA1_L_8 letter3*
 \Rightarrow *SA1_L_8P letter4*
 \Rightarrow *SA1_L_8PP*

Easier Proof that Postcodes are Regular

Can you give an easier proof that the language of postcodes is regular (both left-linear and right-linear)?

Multi-step Regular Grammars

- ▶ In general we can extend regular grammars by allowing productions such as

$$S \longrightarrow abB$$

$$B \longrightarrow aS$$

$$B \longrightarrow baS$$

So instead of having only **one** terminal symbol, we can have **several**.

- ▶ As long as we remain left-linear or right-linear
 - ▶ i.e. the terminal symbols are **always to the right** or **always to the left** of the non-terminal on the right hand side of a rule we obtain grammars which can be reduced to regular grammars.

Lemma II.2.1.2.

Lemma (II.2.1.2.)

1. Assume a grammar G which has only productions of the form

$$A \longrightarrow Bw \text{ or } A \longrightarrow w$$

for some $w \in T^*$, $A, B \in N$. Then $L(G) = L(G')$ for some left-linear grammar G' , which can be computed from G .

2. Assume a grammar G which has only productions of the form

$$A \longrightarrow wB \text{ or } A \longrightarrow w$$

for some $w \in T^*$, $A, B \in N$. Then $L(G) = L(G')$ for some right-linear grammar G' , which can be computed from G .

Multi-step Right-Linear/Left-Linear/Regular Grammars

We call grammars as above multi-step right-linear/left-linear/regular grammars.

Proof Idea for Lemma II.2.1.2.

- ▶ First omit all so called silent productions, i.e. productions of the form $A \rightarrow B$ for some non-terminals A, B .
 - ▶ This requires some work.
- ▶ Then replace in the right-linear case productions

$$A \rightarrow a_1 a_2 \cdots a_n B$$

with $n \geq 2$ by productions

$$\begin{aligned} A &\rightarrow a_1 A_1, \\ A_1 &\rightarrow a_2 A_2, \\ &\dots \\ A_{n-1} &\rightarrow a_n B \end{aligned}$$

for some new nonterminals A_i .

- ▶ Full details can be found in the additional material

Mixing of Left- and Right-Linear

Remark

In a regular grammar we are not allowed to mix left-linear and right-linear grammars. Otherwise we would obtain truly context-free languages.

Example (Mixing Left/Right-Linear Rules)

The following grammar generates the language

$L(G) = ?$

which (as we will see later) is context-free but not regular.

grammar	G
terminals	a, b
nonterminals	S, T
start symbol	S
productions	$S \rightarrow ab$ $S \rightarrow aT$ $T \rightarrow Sb$

II.2.1. Regular Languages (12.2)

II.2.2. Regular Expressions (13.8)

Operators for Forming Languages

Definition

Let $L_1, L_2, L \subseteq T^*$ be languages over the alphabet T .

1. The **concatenation** $L_1.L_2$ of L_1 and L_2 is defined as

$$L_1.L_2 := \{w_1w_2 \mid w_1 \in L_1, w_2 \in L_2\}$$

2. The **union** $L_1 \mid L_2$ of L_1 and L_2 is defined as

$$L_1 \mid L_2 := L_1 \cup L_2$$

The union is sometimes denoted by \pm .

3. The **iteration** or **Kleene-star** L^* of L is defined as

$$L^* := \{w_1w_2 \cdots w_n \mid n \geq 0, w_1, \dots, w_n \in L\}$$

Regular Expressions

Regular expressions are denotations for languages formed from the \emptyset and from the languages $\{a\}$, where a is an element of the alphabet, by using the above mentioned operations.

Regular Expressions

Definition

Let T be an alphabet. We define the set of regular expressions over an alphabet T inductively together with the language $L(E)$ for each regular expression E .

- ▶ \emptyset is a regular expression, $L(\emptyset) := \emptyset$.
- ▶ ϵ is a regular expression, $L(\epsilon) := \{\epsilon\}$.
- ▶ For $a \in T$ we have a is a regular expression, $L(a) := \{a\}$. One usually writes \mathbf{a} for the regular expression, when the symbol is a .
- ▶ If E, F are regular expressions, then
 - ▶ $(E) \mid (F)$ is a regular expression, $L((E) \mid (F)) := L(E) \cup L(F)$.
 - ▶ $(E)(F)$ is a regular expression, $L((E)(F)) = L(E).L(F)$.
 - ▶ $(E)^*$ is a regular expression, $L((E)^*) = L(E)^*$.

We omit unnecessary brackets and usually write $E \mid F$ instead of $(E) \mid (F)$, EF instead of $(E)(F)$, E^* instead of $(E)^*$, if there is no confusion.

Use of Regular Expressions

- ▶ We will usually omit writing $L(E)$, so write

$$(0\ 1)\ 0^*$$

instead of

$$L((0\ 1)\ 0^*)$$

which is

$$(\{0\}.\{1\}).(\{0\})^*$$

which is

$$\{010^n \mid n \in \mathbb{N}\} \text{ or } \{01 \underbrace{0 \cdots 0}_{n \text{ times}} \mid n \in \mathbb{N}\}$$

- ▶ We will as well identify regular expressions which denote the same language. Therefore we can omit more brackets e.g. we can write

$$0\ 1\ 0$$

instead of

$$(0\ 1)\ 0$$

Use of Regular Expressions

- ▶ If the alphabet only contains single characters, we can omit the blank in concatenation, and write

010 instead of **0 1 0**

- ▶ * only refers to the last item, unless there are brackets:
 - ▶ $\mathbf{01}^* = \{0(1^n) \mid n \in \mathbb{N}\}$
 - ▶ $\mathbf{(01)}^* = \{(01)^n \mid n \in \mathbb{N}\}$

Examples of Regular Expressions

- ▶ The set of non-zero digits is defined as

$$\textit{NonzeroDigit} = \mathbf{1} \mid \mathbf{2} \mid \dots \mid \mathbf{9}$$

- ▶ The set of digits is defined as

$$\textit{Digit} = \mathbf{0} \mid \textit{NonZeroDigit}$$

- ▶ The set of numbers without leading zero is

$$\textit{Number} = \mathbf{0} \mid (\textit{NonZeroDigit} \textit{Digit}^*)$$

- ▶ The set of capital letters is defined by

$$\textit{CapitalLetter} = \mathbf{A} \mid \mathbf{B} \mid \dots \mid \mathbf{Z}$$

Examples of Regular Expressions

- ▶ The set of module codes in this department is

$$CSModuleCodes = \mathbf{CS} (- | \mathbf{C} | \mathbf{F} | \mathbf{P}) (\mathbf{0} | \mathbf{1} | \mathbf{2} | \mathbf{3} | \mathbf{M}) \textit{Digit Digit}$$

Regular Expressions are Non-recursive

- ▶ Please note that regular expressions are non-recursive.
For instance in

$$\textit{postcode} = \textit{CapitalLetter CapitalLetter Digit} \sqcup \\ \textit{Digit CapitalLetter CapitalLetter}$$

“postcode” doesn’t occur on the right-hand side.

- ▶ Note that this is different from grammars (including BNF) where recursion is allowed.
For instance we can have productions such as

$$S \longrightarrow aSa$$

or in BNF

$$\langle S \rangle ::= \mathbf{a} \langle S \rangle \mathbf{a}$$

Regular Expressions in Programming

- ▶ Regular Expressions occur very often in programming.
- ▶ They occur in
 - ▶ Linux/Unix (command `grep/egrep`),
 - ▶ in scripting languages (Perl, Python, Ruby),
 - ▶ (one of the main innovations of Ruby over Python was an improved notation `~` for matching of regular expressions),
- ▶ in SQL,
- ▶ are supported in most programming languages by libraries.

Notations for Regular Expressions

- ▶ One writes $[a_1 \cdots a_n]$ for $a_1 \mid \cdots \mid a_n$.
- ▶ One writes $[a - z]$ for $[a, b, c, \dots, z]$ similarly for $[0 - 9]$, $[A - Z]$.
- ▶ $[a - zA - Z] := [a - z][A - Z]$, $[a - z?] := [a - z]?$ etc.
- ▶ One writes L^+ or $L+$ for $L L^*$ (so $L^+ := \{w_1 \cdots w_n \mid n \geq 1, w_1, \dots, w_n \in L\}$, the set of words formed from L by using at least one word in L .
 - ▶ **Question:** Is L^+ the set of non-empty words formed from elements of L^* ?
 - Answer:**
- ▶ Lots of other useful operators for constructing regular expressions have been defined.
- ▶ Each language has its own set and of regular expressions (using often different notations), and its own syntax. Sometimes operators are introduced which go beyond regular languages.

Example Use of Regular Expressions

- ▶ Assume you have files called `automatatheorych1.tex`, `automatatheorych2.tex`, `automatatheorych3.tex`, ...

Concatenation all of them into one file:

```
cssetzer@cs-svr1:> cat automatatheory[0-9].tex >
                    automatatheoryall.tex
```

- ▶ Process lines in a file containing entries separated by “,”, do something if the first field is a student number (a string consisting of digits only). Python code

```
file = open(filename)
regExpStud = re.compile('^ [0-9]*$')
for line in file:
    a = line.split(',')
    if regExpStud.match(a[0]):
        print a[1][:-1] #cut off trailing '\n'
file.close()
```

Example Web Addresses

Consider links in http pages of the form:

```
<a href="http://www.swan.ac.uk/">Swansea University</a>
```

displayed as

[Swansea University](http://www.swan.ac.uk/)

The set of weblinks can be defined as

(in most language `\` would be written as a blank, `"` would be preceded by a `\`, and the whole string would be put into quotation marks):

```
weblinks = <a\ href="[a-zA-Z0-9/.:]">[a-zA-Z0-9\ ]* </a>
```

E.g. in Python one would write

```
weblinks=" <a href=\"[a-zA-Z0-9/.:]*\">[a-zA-Z0-9 ]* </a>"
```

Usage of Regular Expressions in Computer Security

- ▶ In computer security one very often needs to check for occurrences of certain patterns.
- ▶ For instance in order to locate a certain virus, which might consist of 3 pieces of code s_1 , s_2 , s_3 , separated by some normal code, one could search for the regular expression

$$s_1[a - z]^*s_2[a - z]^*s_3$$

(How do you obtain that s_1 , s_2 and s_3 might occur in different order?)

- ▶ Of course in general you need to check for much more sophisticated patterns.

Usage of Regular Expressions in Computer Security

- ▶ In order to check that a password is safe enough, which might mean it consists of digits and lower case characters, and at least one digit and one lower case character, would mean that you whether it matches

$$(([a-z] | [0-9])^* [a-z] ([a-z] | [0-9])^* [0-9] ([a-z] | [0-9])^*) | (([a-z] | [0-9])^* [0-9] ([a-z] | [0-9])^* [a-z] ([a-z] | [0-9])^*)$$

Of course you would usually use a much more sophisticated regular expression.

Usage of Regular Expressions in Computer Security

- ▶ Detecting in request certain malicious patterns in requests from the outside can often be expressed as a regular expression and you search for matches in this income stream which match that expression.

SQL Injection

- ▶ Regular expression can be used to detect attempts of SQL injection.
- ▶ Example of SQL injection (from Wikipedia on SQL Injection):
Assume the following statement in a code

statement =

```
“SELECT * FROM users WHERE username = ’” + username +  
    “ ’ and password = ’” + password + “ ’;”
```

- ▶ This statement is supposed to be sent to the SQL server.
Then one checks the resulting entries for whether the supplied password matches modulo encryption one of the password entries for that user name.

SQL Injection

- ▶ Assume an attacker tries to login with password

' or '1'='1

Then the statement sent to the SQL server will be

```
SELECT * FROM users WHERE username = 'username'  
and passport = " or '1'='1';
```

which matches all users with username.

SQL Injection

- ▶ This might allow you to check whether your password matches any user, which makes it more likely to get a match and allow you to login.
- ▶ In order to avoid such kind of attack you can check whether the username matches any malicious pattern.
Such patterns can be expressed by regular expressions.

Usage of Regular Expressions in Computer Security

- ▶ The above were just some (here very simple) examples how regular expressions can be used to detect in computer security certain patterns corresponding to attacks or weaknesses of a system.

Closure of Regular Languages

The main lemma for showing that regular expressions define regular languages is as follows:

Lemma (II.2.2.1.)

Let G, G' be both left-linear grammars or both right-linear grammars. Then we can define a left-linear or right-linear grammars G_i s.t.

1. $L(G_1) = L(G) \mid L(G')$,
2. $L(G_2) = L(G).L(G')$,
3. $L(G_3) = L(G)^*$.

These grammars can be computed from G and G' .

Proof

A proof can be found in the additional material for this subsection.

Regular Expressions define Regular Languages

Lemma (II.2.2.2.)

Let E be a regular Expression. Then there exist both left-linear and right-linear grammars G, G' s.t.

$$L(E) = L(G) = L(G')$$

G and G' can be computed from L .

Proof: By Lemma II.2.2.1, and the fact that the finite languages $\emptyset, \{\epsilon\}$ and $\{a\}$ are regular.

Full details can be found in Additional Material.