

ECOTE – final project

Semester: 21L

Author: Wiktor Bara

Subject: NFA construction using Thompson algorithm

I. General overview and assumptions

The task is to design a program reading regular expression, then constructing Nondeterministic Finite Automata, and checking if input strings are generated by this expression using already created automata. The program should take two files as an input, first one with one regular expression and second one with input strings, and return one file as an output that contains information whether given strings could be generated by given regular expression.

I assumed that the alphabet for the regular expression contains all numbers and letters (both small and uppercase). The end of line ('\n') will be treated like end of the expression, and any other characters, that are not contained in regular expression definition, will be treated as an error.

II. Functional requirements

The general purpose for NFA constructor is to create automata and check whether the strings that are given can be created by this automata. This means that the class NFA is needed, in which based on the regular expression the automata is stored to be used multiple times. This data structure is further described in *Data structures* section.

The NFA constructor evaluates each character from the regular expression file, changing its behavior whenever it encounters special symbol. The special symbols of the NFA constructor are:

Symbol	Meaning	Priority
(left parenthesis that opens block with greater calculation priority	-
)	right parenthesis that closes the block with greater calculation priority	-
?	zero or one characters	5
+	one or more characters	4
*	zero or more characters	3
.	concatenation	2
	alteration	1

The symbols have their priorities that force to generate NFA in the appropriate way and determines the order of the characters.

The NFA creator can also output the errors and warnings, when incorrect, unexpected or inadvisable input is encountered. The program counts the lines in the input file to make the errors more useful and helpful for the user. The errors and warnings will be described in the *Input/output description* section.

III. Implementation

General architecture

The program consists of five main modules:

- Main Module – responsible for the Command Line Interface and File I/O
- NFA Creator Module – responsible for creation of the automata
- String Comparison Module – responsible for checking whether given strings can be created by automata
- Error Module – responsible for the list of descriptions of errors and warnings

The class diagram and description are presented below:

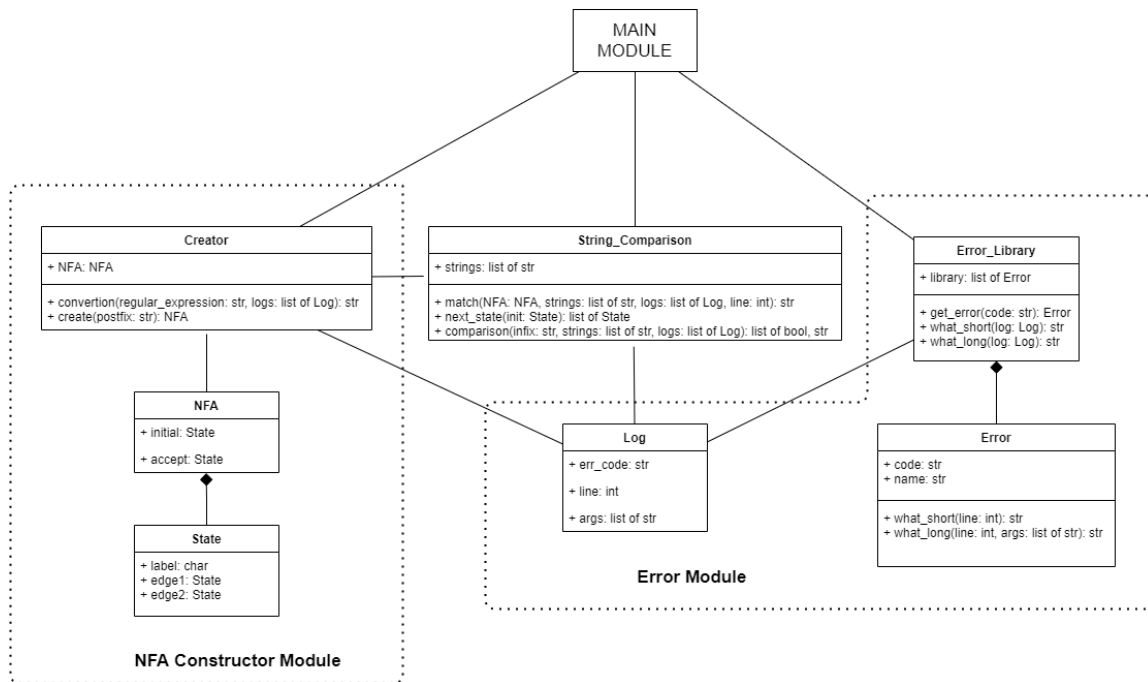


Figure 1. Class Diagram

Class Creator

NFA: NFA	NFA object
convert(regular_expression: str, logs: list of Log): str	The function transform source text from file with regular expression and converts it to the Postfix Notation (Reverse Polish notation) using Shunting-Yard Algorithm.
compile(postfix: str): NFA	The function takes already converted to Postfix Notation regular expression and creates NFA object that corresponds to the input

Class NFA

initial: State	The initial state of NFA that can be connected to another NFA's, it is the first state of the NFA
accept: State	The accept state that can be connected to the other NFA's, it is the last state of the NFA

Class State

label: str	label defines the argument that has to be given to proceed the transition, when null or None the transition is equivalent to ϵ transition
edge1: State	edge1 defines one (first) of the connections to other states
edge2: State	edge2 defines one (second) of the connections to other states

Class String_Comparison

strings: list of str	The strings that was in the input file and where separated by end of line (' $\backslash n$ ') are now in the list and every one of them will be checked whether can be created by the NFA
match(NFA: NFA, strings: list of str, logs: list of Log, line: int): str	Function checks whether the strings can be created by the NFA
next_state(init: State): list of State	returns set of States that can be reached from the given state with ϵ transition
comparison(infix: str, strings: list of str, logs: list of Log): list of bool, str	returns the list of boolean values created based on checking whether given strings can be generated by RE

Class Log

err_code: str	The error/warning code encountered
line: int	The line at which the error/warning occurred
args: list of str	List of additional arguments, used to produce the verbose error description

Class Error_Library

library: list of Error	A list of Error objects
get_error(code: str): Error	A method returning an Error object given the error code
what_short(log: Log): str	A method generating a short error description based on a Log object
what_long(log: Log): str	A method generating a verbose error description based on a Log object

Class Error

code: str	The error code
name: str	The name of the error
what_short(log: Log): str	Method returning short error description
what_long(log: Log): str	Method returning verbose error description

Data structures

Structure	Description
ErrorLibrary	List of error and warning definitions. It consists of Error objects. As the error codes are unique, it should be a list that allows such behavior
NFA	The data structure that consists of two States initial (the first state) and accept (the last state). The States that are between those two given can be accessed by the direct connections between States as an objects.
State	The data structure that connects to other instances of itself. It corresponds to nodes and edges in the NFA diagram.

Module descriptions

The most important part of the modules are presented below: the execution of the main module, the conversion to Postfix Notation, creation of the NFA

Main Module:

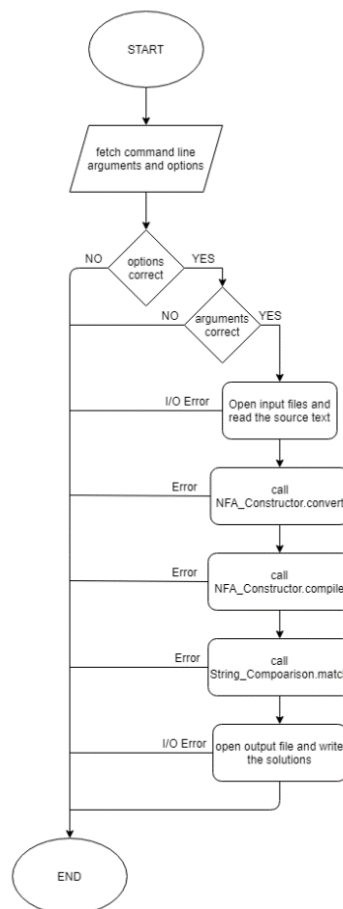


Figure 2. Main module

Conversion to Postfix Notation:

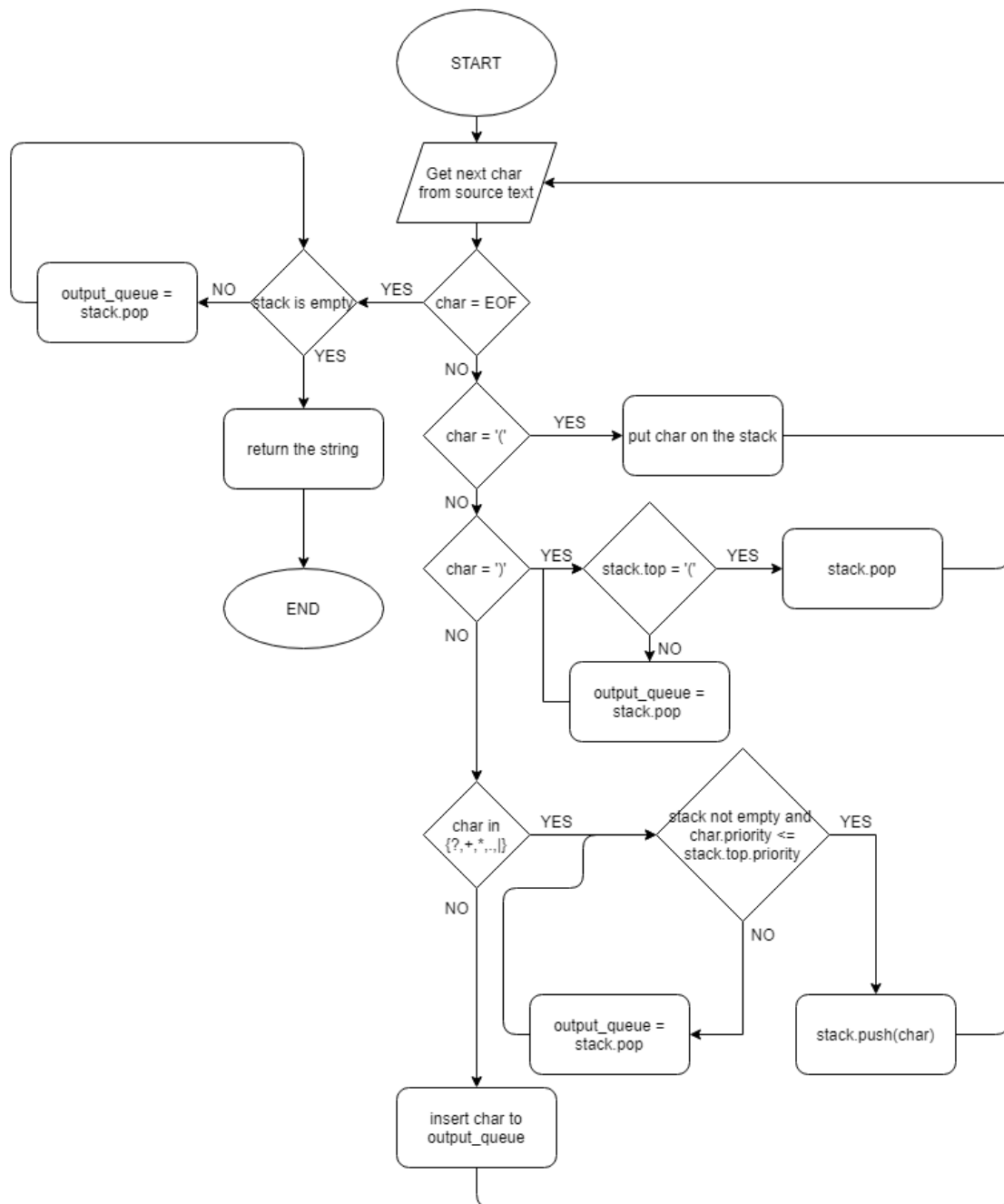


Figure 3. Conversion to Postfix Notation

Creation of the NFA:

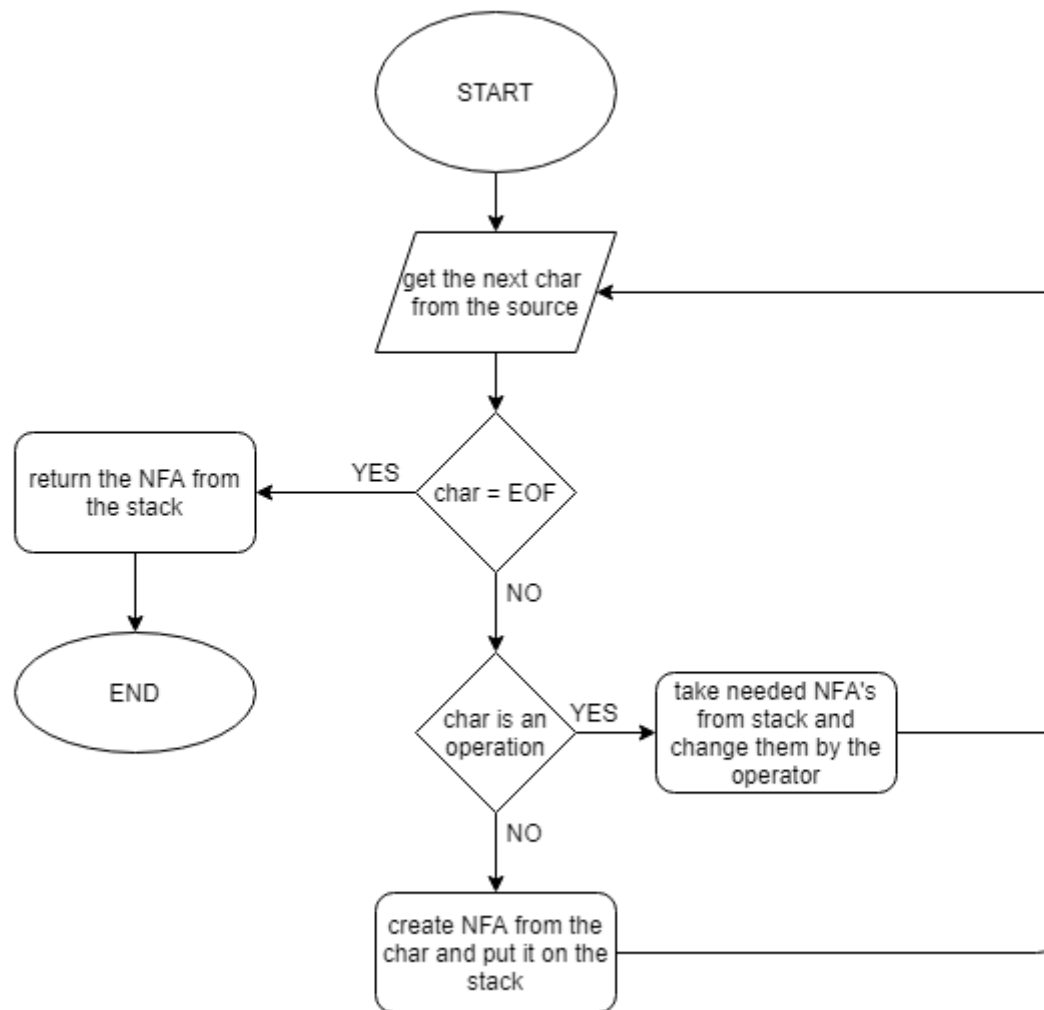


Figure 4. Creation of the NFA

The flowcharts for conversion and NFA creation does not include error handling. At any point in the execution, when an error occurs, an exception will be thrown, or a warning will be added to the log list.

Input/output description

The program uses two text files as input, and outputs text into another file. The NFA creation and string checking works as previously described. Below are some basic examples of usage of this program.

Input	Output
file1: 1*0	False – RE: 1*0, String: 1111
file2: 1111	True – RE: 1*0, String: 10
10	True – RE: 1*0, String: 0
0	False – RE: 1*0, String: 100
100	

Example 1: Basic check of some strings

Input	Output
file1: 1*0 file2: \$22 1110	Error e30: Not defined character in String file

Example 2: Incorrect usage – the string characters not in alphabet (*more about errors in Errors and Warnings section*)

Input	Output
file1: 1)*0 file2: 22 1110	Error e11: Lack of left parenthesis

Example 3: Incorrect usage – lack of left parenthesis (*more about errors in Errors and Warnings section*)

Command Line Interface

The program is operated with a command line interface (CLI). The basic usage of CLI is presented below:

```
> const regural_expression_file strings_file output_file
```

The program takes 'regural_expression_file' as an input to the NFA constructor and 'strings_file' to the strings checker. The output text is written into 'output_file'. The errors and warnings are outputted to stdout. When there are only two arguments the program will write the outputs to the default file named 'output' while when there is less than two or too many arguments the error is raised. The program also supports options to control the error and warning output, that can be inserted before the I/O files. The option summary can be seen in the table below.

Option	Short version	Arguments	Description
--silent	-s	None	Mutes the error and warning output
--nowarn	-n	None	Mutes warning output
--output	-o	Filename	Redirects the error and warning to a file
--verbose	-v	None	Makes the error and warning output more verbose to give a clearer error information

Choosing a non-existent option or incompatible options (like -s + -v) will result in an error and the program will abort execution.

Errors and Warnings

During computation the program can encounter errors and warnings that are caused by incorrect or unexpected input. In case when there is an error raised the program informs user of the type of error and stops further execution. In case of a warning program informs user of the type of warning and continues execution. The errors and warnings are presented in the tables below.

Category	Code	Error Name	Description
<i>Conversion</i>	e10	Not defined character in RE	The regular expression contains character that is not number or a letter or one of the special symbols { (,) , ? , + , * , . , }
	e11	Lack of left parenthesis	The regular expression reader encountered ')' while there was not corresponding symbol '(' before
	e12	Lack of right parenthesis	The regular expression reader encountered '(' that was not closed by the ')' symbol
	e13	Invalid operation symbol placement - start	The regular expression reader encounter on one on the following symbols at the beginning of the RE { ? , + , * , . , }
	e14	Invalid operation symbol placement - neighbour	Pair of the following symbols were encountered next to each { . , } other
	e15	Invalid operation symbol placement - pharenthesis	One of the following was placed after '(' symbol { ? , + , * , . , } or one of the following was placed before ')' symbol { . , }
	e16	Invalid operation symbol placement – neighbour2	The repetition symbol { ? , + , * }, was encountered after non-repetition symbol { . , }.
	e17	Repetition operators close to each other	Pair of the following symbols were encountered next to each other what can cause unexpected results, { ? , + , * } and deadlock.
<i>Strings</i>	e30	Not defined character in String file	The strings contains character that is not a number or a letter
<i>Other</i>	e98	I/O Error	There was an error with the file I/O

Category	Code	Warning Name	Description
<i>Conversion</i>	w10	RE Whitespace	The regular expression contains whitespace character at the end, anything that is placed after that character will not be read
<i>String</i>	w30	Whitespace argument	One of the Strings is a whitespace
<i>CLI</i>	w80	Overwrite Warning	The input file is the same as the output file

Others

I intend to implement project in Python.

IV. Functional test cases

The following test cases check all the warnings and errors available in the MacroGenerator and a few correct usage cases. The tests don't include the CLI tests.

Correct usage

Test case	Input	Output
Basic RE with one string that can be created	file1: 1* file2: 111	True – RE: 1*, String: 111
Basic RE with one string that cannot be created	file1: 1* file2: 110	False – RE: 1*, String: 110
Basic RE with multiple string that can be created	file1: 1* file2: 111 1	True – RE: 1*, String: 111 True – RE: 1*, String: 1
Basic RE with multiple string that cannot be created	file1: 1* file2: 110 0xa	False – RE: 1*, String: 110 False – RE: 1*, String: 0xa
Basic RE with multiple string	file1: 1*.(0 a)*.x file2: 110ax aaaax 1111xx	True – RE: 1*.(0 a)+.x , String: 110ax True – RE: 1*.(0 a)+.x , String: aaaax False – RE: 1*.(0 a)+.x , String: 1111xx

Errors

Test case	Input	Output
Not defined character in RE	file1: 1*/a file2: 111	<i>error e10</i>
Lack of left parenthesis	file1: 1*a c) file2: 110	<i>error e11</i>
Lack of right parenthesis	file1: 1*(a c file2: 111	<i>error e12</i>
Invalid operation symbol placement - start	file1: +1* file2: 111	<i>error e13</i>
Invalid operation symbol placement - neighbour	file1: 1*a. c file2: 111	<i>error e14</i>
Invalid operation symbol placement - parenthesis	file1: 1(*a.)c file2: 111	<i>error e15</i>
Invalid operation symbol placement – neighbour2	file1: 1.*a file2: 110	<i>error e16</i>
Repetition operators close to each other	file1: 1*+ file2: 111	<i>error e17</i>
Not defined character in String file	file1: 1* file2: 11*1	<i>error e30</i>

Warnings

Test case	Input	Output
RE Whitespace	file1: 1* xa file2: 111	True – RE: 1*, String: 111 <i>warning w10</i>
Whitespace argument (I present '_' as a ' ')	file1: 1* file2: _	True – RE: 1*, String: _ <i>warning w30</i>

V. Python Implementation

The program was implemented using **Python**, version **3.8.7**, and ran and tested under Windows 10 environment. The creation of NFA and string comparison part uses standard functionalities of Python. The main module, responsible for the CLI and File I/O uses the Python libraries `sys` and `optparse`. For testing the library `unittest` was used. The NFA creator and error modules can be imported and used outside from the supplied CLI application. The source code is available in the `src` directory, inside which are `main.py`, the main CLI application, `test.py`, which run the tests, and separate folders with `Errors`, `Creator` and `StringComparison` modules. The implementation is available in a repository on github:

https://github.com/lewiktorbara/ECOTE_NFA_creation

File Structure

```
└── src
    ├── main.py
    ├── test.py
    ├── re.txt
    ├── err_re.txt
    ├── strings.txt
    ├── Creator
    │   ├── __init__.py
    │   └── creator.py
    ├── Errors
    │   ├── __init__.py
    │   ├── error_library.py
    │   ├── error.py
    │   └── log.py
    ├── StringComparison
    │   ├── __init__.py
    │   └── comparator.py
    └── Tests
        ├── __init__.py
        └── nfa_test.py
```

Running the program

The tests can be run by running the `src/test.py` file

```
D:\Database\Studia\Sem8\ECOTE\Project>python src/test.py
```

.....

```
Ran 16 tests in 0.002s
```

OK

Figure 5. Running the unit tests

The main program is run by running `src/main.py` file, providing input files as arguments. A text files `re.txt`, `err_re.txt` and `strings.txt` are provided, which can be used to run the program. Detailed description of the CLI is available in the *Command Line Interface* section.

```
D:\Database\Studia\Sem8\ECOTE\Project\src>python main.py re.txt strings.txt
```

```
True - RE: 1.a+.(b|c)*, String: 1aaaaaaaa
```

```
True - RE: 1.a+.(b|c)*, String: 1abcbcbcbcbc x
```

```
False - RE: 1.a+.(b|c)*, String: 1bcbcbcb
```

```
False - RE: 1.a+.(b|c)*, String: 1bbbbbbbbbbcccccccccccc
```

```
Execution completed with 1 warning:
```

```
w30 Whitespace argument at line 1.
```

Figure 6. Example of running the program with correct input

```
D:\Database\Studia\Sem8\ECOTE\Project\src>python main.py err_re.txt strings.txt
```

Execution unsuccessful.

```
e12 Lack of right parenthesis at line 0.
```

Figure 7. Example of running the program with incorrect input