

Laboratory Work Nr.4 Regular Expressions

Cretu Cristian

April 9, 2024

Theory

Regular expressions, often abbreviated as regex or regexp, are powerful tools used in computer science and programming for pattern matching within strings of text. They provide a concise and flexible means of searching, extracting, and manipulating textual data based on specific patterns. Utilized across various programming languages and text processing utilities, regular expressions enable developers to perform tasks such as validation of input, searching for specific patterns within large datasets, and text manipulation with precision and efficiency. The syntax of regular expressions consists of a combination of literal characters and metacharacters, forming patterns that define the desired matches. With their versatility and widespread adoption, regular expressions serve as an indispensable tool for tasks ranging from simple string manipulation to complex data extraction and transformation. However, mastering regular expressions requires understanding their syntax, metacharacters, and application-specific nuances to leverage their full potential effectively.

Objectives

1. Write and cover what regular expressions are, what they are used for.
2. Take your variant code
3. Write a code that will generate valid combinations of symbols conform given regular expressions (examples will be shown).
4. In case you have an example, where symbol may be written undefined number of times, take a limit of 5 times (to evade generation of extremely long combinations).
5. Bonus point: write a function that will show sequence of processing regular expression (like, what you do first, second and so on)

Implementation description

RegEx Operators

I have created custom functions for handling different Regex operators, particularly the +, |, *, ^ operators.

For example. here is the OR operator:

```
1 def or_operator(options: list[str] = None):
2     symbol = random.choice(options)
3     print(f'Choose {symbol}')
4     return symbol
```

The OR operator would perform an operation on a group of symbols and randomly choose one symbol from a list of those. For example, giving it a list of ['P','Q','R'], it will choose one of those. In a Regular Expression, this would be noted as (P—Q—R), so you choose one of them.

Here is the Star Operator

```
1 def star_operator(symbol: str):
2     times = random.randint(0, infity)
3     answer = symbol*times
4     print(f'Repeat {symbol} {times} times')
5     return answer
```

The * - STAR operator means that a symbol can be repeated 0 or an infinite times. In this laboratory work, the “infity” variable is a global one and is set to 5, so that we can represent the Regular Expression correctly and will not get too many repeating symbols.

Other operators are:

- + - symbol is allowed to be repeated 1 or an infinity amount of times.
- ^x - symbol should be repeated x times consecutively.
- ? - symbol may or may not appear. My implementation just randomly assigns a truth value that represents whether or not the symbol is included or not.

After defining these functions, we have to somehow parse the Regex notation to be able to perform these above mentioned operations on the rules specified. To do that, we search the notation for parenthesis so that we can find groups of symbols to choose from, and other symbols are represented as simple strings.

```
1 def compute_regex(parsed_regex: list[str] = None):
2     if parsed_regex is None:
3         return None
4     indices_to_skip = []
5     answer = ""
6     skipped_simbols = ["(", ")", "*", "+", "?", "^"]
7     for i in range(len(parsed_regex)):
8         if i in indices_to_skip:
```

```

9         continue
10        part = parsed_regex[i]
11        if part not in skipped_symbols:
12            if "|" in part:

```

Here we check if this part is a group of symbols to choose from, and if so, we randomly choose one of those

```

1            current_token = or_operator(part.split("|"))
2        else:

```

Otherwise, we just take the current symbol

```

1            current_token = part
2            print(f"Current Token: {current_token}")
3            if i+1 < len(parsed_regex):

```

Then we look which operation follows the current token

```

1            next_part = parsed_regex[i+1]
2            if next_part == "*":
3                answer += star_operator(current_token)
4            elif next_part == "+":
5                answer += plus_operator(current_token)
6            elif next_part == "?":
7                answer += question_operator(
            current_token)
8            elif next_part == "^":

```

We will skip later the symbol that follows the ^ operator, since we don't want to process it as a separate symbol.

```

1                power = int(parsed_regex[i+2])
2                answer += n_operator(current_token,
3                power)
4                indices_to_skip.append(i+2)
5            else:
6                answer += current_token
7        else:
8            answer += current_token
9        return answer

```

After that, we can proceed to traverse the parsed regex and perform the operations step by step to obtain the desired string. This function will take a part of the parsed Regex by one, and perform the necessary operations on it, and append them one over one to the final string obtained by the Regex.

```

1 def compute_regex(parsed_regex: list[str] = None):
2     if parsed_regex is None:
3         return None
4     indices_to_skip = []
5     answer = ""
6     skipped_symbols = ["(", ")", "*", "+", "?", "^"]
7     for i in range(len(parsed_regex)):

```

```

8     if i in indices_to_skip:
9         continue
10    part = parsed_regex[i]
11    if part not in skipped_symbols:
12        if "|" in part:
13            current_token = or_operator(part.split("|"))
14        else:
15            current_token = part
16            print(f"Current Token: {current_token}")
17            if i+1 < len(parsed_regex):
18                next_part = parsed_regex[i+1]
19                if next_part == "*":
20                    answer += star_operator(current_token)
21                elif next_part == "+":
22                    answer += plus_operator(current_token)
23                elif next_part == "?":
24                    answer += question_operator(
25                        current_token)
26                elif next_part == "^":
27                    power = int(parsed_regex[i+2])
28                    answer += n_operator(current_token,
29                        power)
30                    indices_to_skip.append(i+2)
31                else:
32                    answer += current_token
33            else:
34                answer += current_token
35    return answer

```

This function makes the difference between a group of symbols to choose from and simple symbols, and looks at the next symbol in list to see if there is some operation that is needed to perform. Also, for the power operator, it will select as power the part next to the next part and include it into a list of indices to avoid, so that it would not pass over that symbol again and not append it as a single letter.

Execution

The first Regex's execution would look like this

```

1  Regex: 0(P|Q|R)+2(3|4)
2  Parsed Regex: ['0', 'P|Q|R', '+', '2', '3|4']
3  Current Token: 0
4  Choose Q from ['P', 'Q', 'R']
5  Repeat Q 2 times
6  Current Token: 2
7  Choose 3 from ['3', '4']
8  Computed Regex: 0QQ23

```

Conclusions

In conclusion, this laboratory work provided valuable insights into regular expressions, emphasizing their role in pattern matching within text. Practical implementation involved creating custom functions for handling regex operators and generating valid symbol combinations. Exploring operators like `*`, `+`, `—`, `?`, and `^` enhanced understanding of pattern definition and repetition control. Future enhancements could include supporting additional regex operators, optimizing code for efficiency, and implementing error handling mechanisms. Overall, this experience deepened comprehension of regex principles and their practical application, setting a foundation for further exploration in future projects.