

Laboratory Work No. 4

Course: Formal Languages & Finite Automata

Student: Iordan Liviu, Group: FAF-223

Topic: Regular expressions

Abstract

This laboratory work explored the use of regular expressions for text pattern matching. A Python program was developed to generate valid combinations of symbols based on a complex assigned regular expression. The code successfully produced the expected combinations, and a function was implemented to provide a step-by-step explanation of the regular expression processing logic.

Introduction

Regular expressions (often shortened to "regex") are powerful tools within computer science, offering a concise and flexible way to define patterns within text. They are used extensively in tasks like:

- Text Search and Validation: Locating specific patterns within large bodies of text and ensuring input data (e.g., email addresses, phone numbers) conforms to specific rules.
- Code Refactoring and Text Manipulation: Finding and replacing text patterns in programming code or complex documents.

- **Data Extraction:** Pulling out specific information from unstructured text for further processing and analysis.

Regular expressions have a formal foundation in the theory of formal languages and finite automata. A regular expression defines a regular language, which can be recognized by a finite automaton. This theoretical basis underlies the efficiency and expressiveness of regular expressions for practical pattern-matching tasks.

This laboratory work investigates the practical application of regular expressions. Through code implementation and analysis, we will explore how regular expressions can be constructed to match specific patterns, generate valid combinations, and understand the step-by-step logic involved in their processing.

Key Concepts

To fully understand this work, it's useful to be familiar with the following concepts:

- **Metacharacters:** Special characters in regular expressions (e.g., `*`, `+`, `?`, `|`) that have specific pattern matching meanings.
- **Character Classes:** Sets of characters within square brackets (e.g., `[a-z]`, `[0-9]`) that match a single character from the given set.
- **Quantifiers:** Symbols that specify the repetition of pattern elements (e.g., `*` zero or more times, `+` one or more times).
- **Alternation:** Using the pipe symbol (`|`) to offer options within a pattern.

Objectives:

Write and cover what regular expressions are, what they are used for; Below you will find 3 complex regular expressions per each variant. Take a variant depending on your number in the list of students and do the following:

- a. Write a code that will generate valid combinations of symbols conforms to given regular expressions (examples will be shown).
- b. 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);
- c. Bonus point: write a function that will show sequence of processing regular expression (like, what you do first, second and so on)

Regular Expression Variant explanation

- `M?`: Matches zero or one occurrence of the letter 'M'.
- `N^2`: Matches exactly two occurrences of the letter 'N'.
- `(O|P)^3`: Matches either 'O' or 'P' repeated three times (e.g., "OOO", "PPP", "OPO").
- `*Q^**`: Matches zero or more occurrences of the letter 'Q'.
- `R^+`: Matches one or more occurrences of the letter 'R'.
- `(X|Y|Z)^3`: Matches any sequence of exactly three characters from the set {X, Y, Z} (e.g., "XYZ", "ZYX", "YYX").

- `8^+`: Matches one or more occurrences of the digit '8'.
- `(9|0)`: Matches either the digit '9' or the digit '0'.
- `(H|I)`: Matches either the letter 'H' or the letter 'I'.
- `(J|K)`: Matches either the letter 'J' or the letter 'K'.
- `L*N?`: Matches zero or more occurrences of the letter 'L' followed by zero or one occurrence of the letter 'N'.

Implementation

- **Function `generate_string`**

```
def generate_string(pattern: str) -> str:
    result = ""
    i = 0
    while i < len(pattern):
        char = pattern[i]
        if char == '(':
            j = i + 1
            subpattern = ""
            while pattern[j] != ')':
                subpattern += pattern[j]
                j += 1
            choices = subpattern.split('|')
            subpattern = random.choice(choices)
            repeat = 1
            i = j + 1
            # Handling special characters '^', '*', '+', '?'
            if j + 1 < len(pattern):
                next_char = pattern[j + 1]
                if next_char == '^':
                    if pattern[j + 2].isdigit():
                        repeat = int(pattern[j + 2])
                        i = j + 3
                    elif next_char == '*':
                        repeat = random.randint(0, 3)
                        i = j + 2
                    elif next_char == '+':
                        repeat = random.randint(1, 3)
```

```

        i = j + 2
    elif next_char == '?':
        repeat = random.randint(0, 1)
        i = j + 2
    result += subpattern * repeat
    continue
# Handling special characters '*', '+', '?', '^'
next_char = pattern[i + 1] if i + 1 < len(pattern) else ""
if next_char == '^':
    result += char * int(pattern[i+2])
    i += 3
    continue
elif next_char == '*':
    result += char * random.randint(0, 5)
    i += 2
    continue
elif next_char == '+':
    result += char * random.randint(1, 5)
    i += 2
    continue
elif next_char == '?':
    result += char * random.randint(0, 1)
    i += 2
    continue
result += char
i += 1
return result

```

The `generate_string` function takes a pattern string as input and constructs a new string based on that pattern. It iterates through each character of the pattern using a while loop. If it encounters an opening parenthesis '(', it collects characters until the closing parenthesis ')' to form a subpattern with multiple choices separated by '|'. It then selects one of these choices randomly and handles potential repetition specified by special characters ('^', '*', '+', '?') that immediately follow the closing parenthesis. For other characters, it checks for similar special characters that denote repetition and appends the appropriate number of repetitions (chosen randomly within specified ranges) to the result string. The function continues this process until it has processed the entire pattern and returns the generated string. The goal is to dynamically construct varied strings based on the provided pattern and its specified rules for character choices and repetition.

- **Function `explain_pattern_processing`**

```

def explain_pattern_processing(pattern: str) -> list:
    explanation = []
    i = 0

```

```

while i < len(pattern):
    char = pattern[i]
    if char == '(':
        j = i + 1
        subpattern = ""
        while pattern[j] != ')':
            subpattern += pattern[j]
            j += 1
        choices = subpattern.split('|')
        explanation.append(f"Select one of: {choices}")
        i = j + 1
        # Handling special characters '^', '*', '+', '?'
        if j + 1 < len(pattern):
            next_char = pattern[j + 1]
            if next_char == '^':
                repeat = int(pattern[j + 2])
                explanation.append(f"Repeat the selection {repeat} times")
                i = j + 3
            elif next_char == '*':
                explanation.append("Randomly repeat zero to three times")
                i = j + 2
            elif next_char == '+':
                explanation.append("Randomly repeat one to three times")
                i = j + 2
            elif next_char == '?':
                explanation.append("Randomly repeat zero or one time")
                i = j + 2
            continue
        # Handling special characters '*', '+', '?', '^'
        next_char = pattern[i + 1] if i + 1 < len(pattern) else ""
        if next_char == '^':
            repeat = int(pattern[i + 2])
            explanation.append(f"Repeat {char} {repeat} times")
            i += 3
            continue
        elif next_char == '*':
            explanation.append(f"Randomly repeat {char} zero to five times")
            i += 2
            continue
        elif next_char == '+':
            explanation.append(f"Randomly repeat {char} one to five times")
            i += 2
            continue
        elif next_char == '?':
            explanation.append(f"Randomly repeat {char} zero or one time")
            i += 2
            continue
        explanation.append(f"Use character: {char}")
        i += 1

return explanation

```

The `explain_pattern_processing` function analyzes a pattern string character by character to provide a step-by-step explanation of how the pattern is processed. It initializes an empty list `explanation` to store each explanation string. The function iterates through the pattern using a while loop until it reaches the end of the pattern string. For each character `char` in the pattern, it checks if `char` is an opening parenthesis '('. If so, it collects characters until the closing parenthesis ')' to form a subpattern containing choices separated by '|'. The function then appends an explanation string indicating the available choices within the subpattern. It also handles special characters ('^', '*', '+', '?') that immediately follow the closing parenthesis by appending explanations for the specified repetitions. If `char` is not an opening parenthesis, the function checks the next character (`next_char`) to determine if it indicates a repetition. If so, it appends an explanation string for the random repetition behavior specified by the special character. If neither of these conditions is met, the function appends an explanation string indicating the use of the current character `char`. The function continues this process until it has processed the entire pattern string and then returns the list of explanations detailing each step taken during the pattern processing. This function is designed to provide a clear breakdown of how the input pattern is interpreted and transformed into a series of operations to generate a string.

- **Verifying the code and displaying the steps (shown for pattern one, but it is the same for other patterns, just change the variable names)**

```
# Pattern 1

print("Explanation for Pattern 1:")

explanations = explain_pattern_processing(pattern1)

for step in explanations:

    print(step)

print("\nExamples of words made after pattern 1:")
```

```
generated_strings = [generate_string(pattern1) for _ in range(5)]  
  
print(", ".join(generated_strings))
```

The provided code snippet showcases the functionality of two custom functions, `explain_pattern_processing` and `generate_string`, applied to a specific pattern (`pattern1`). Firstly, the code calls `explain_pattern_processing(pattern1)` to generate a step-by-step explanation of how `pattern1` is processed, storing the explanations in a list named `explanations`. It then iterates through `explanations` to print out each step, elucidating the interpretation and processing logic of `pattern1`. Subsequently, the code utilizes the `generate_string` function within a list comprehension `[generate_string(pattern1) for _ in range(5)]` to generate five example words based on `pattern1`. These generated words are stored in `generated_strings`, which are subsequently printed to demonstrate the diversity of outcomes achievable with `pattern1` using the `generate_string` function. Together, this sequence of actions illustrates the interpretive and generative capabilities of the provided custom functions within the context of a specific pattern, offering insights into their underlying mechanisms and utility for text generation tasks.

Difficulties and Solutions

Difficulties:

- Understanding complex regular expressions: Complex regular expressions with nested parentheses and various special characters can be challenging to interpret.
- Handling undefined repetition: Regular expressions can allow characters to be repeated a certain number of times (e.g., `x*`). Implementing logic to limit repetitions within the code can be tricky.

- Debugging matching issues: Ensuring the code accurately matches strings based on the regular expression can involve debugging and testing with various inputs.

Solutions:

- Break down complex expressions: Analyze the regular expression step-by-step, identifying patterns and their meanings. Utilize online resources or tools for visualizing regular expressions.
- Implement limits for repetitions: When encountering * or +, modify the code to create a loop that iterates a specific number of times (based on the requirement, like limiting to 5 repetitions in this case).
- Test with various inputs: Create a variety of test cases with valid and invalid strings to ensure the code correctly matches intended patterns and rejects unexpected inputs.

Conclusions

In this laboratory work, we delved into the realm of regular expressions, a powerful tool for pattern matching and text processing. Our objectives were to understand what regular expressions are, their applications, and to implement a Python code capable of generating valid combinations of symbols based on complex regular expressions.

We successfully achieved these objectives by developing a Python script that utilizes regular expression parsing and itertools module to generate valid combinations. Additionally, we went a step further by implementing a function to provide a detailed sequence of processing steps for each regular expression, enhancing our understanding of how the patterns are interpreted.

Throughout the project, we encountered several challenges, primarily in understanding and parsing complex regular expressions. However, through collaborative efforts and perseverance, we overcome these hurdles, leading to a robust implementation.

By adhering to the evaluation criteria, we ensured the project's accessibility and transparency. The codebase is hosted in a public GitHub repository, enabling easy access for evaluation and future reference. Detailed explanations of the performed work, including the functionality of the code and the sequence of processing steps, were provided to facilitate understanding.

In conclusion, this laboratory work not only deepened our understanding of regular expressions but also honed our programming skills in Python. It provided hands-on experience in tackling real-world text processing tasks, which will undoubtedly prove invaluable in our academic and professional endeavors.