

Natural Language Processing

An Introduction

Abdelbacet Mhamdi

2025-01-23

ISET Bizerte

1. Introduction to Regular Expressions (Regex)

1. Introduction to Regular Expressions (Regex)

1.1 What are Regular Expressions?

Regular expressions are powerful patterns used to match, search, and manipulate text strings. They provide a standardized way to describe search patterns in text, making them an essential tool in programming, text processing, and data validation.

1.2 Core Concepts

1.2.1 Pattern Matching

A regex pattern is a sequence of characters that defines a search pattern. These patterns can be:

- Literal characters that match themselves;
- Special characters (metacharacters) with special meanings;
- Combinations of both.

1.2 Core Concepts

1.2.2 Basic Metacharacters

Metacharacter	Description	Example
.	Matches any character except newline	a.c matches "abc", "a1c", "a@c"
^	Matches start of string	^Hello matches "Hello World"
\$	Matches end of string	world\$ matches "Hello world"
*	Matches 0 or more occurrences	ab*c matches "ac", "abc", "abbc"
+	Matches 1 or more occurrences	ab+c matches "abc", "abbc" but not "ac"
?	Matches 0 or 1 occurrence	ab?c matches "ac" and "abc"
\	Escapes special characters	\. matches literal dot

1.3 Common Use Cases

1. Search Operations

- Advanced find/replace operations;
- Pattern matching in large text files;
- Content filtering.

2. Text Processing

- Finding patterns in text;
- Replacing specific text patterns;
- Extracting information;
- Parsing log files.


3. Data Validation

- Email addresses;
- Phone numbers;
- Postal codes;
- Passwords;
- URLs.

1.4 Advanced Concepts

1.4.1 Character Classes

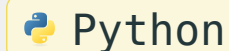
```
1 # Character class examples
2 pattern = r'[aeiou]' # Matches any vowel
3 pattern = r'[0-9]'   # Matches any digit
4 pattern = r'^0-9]'   # Matches any non-digit
```

 Python

1.4 Advanced Concepts

1.4.2 Quantifiers and Groups

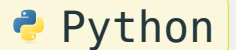
```
1 # Quantifiers
2 pattern = r'\d{3}'      # Exactly 3 digits
3 pattern = r'\d{2,4}'    # Between 2 and 4 digits
4 pattern = r'\d{2,}'     # 2 or more digits
5
6 # Groups
7 pattern = r'(\w+)\s+\1'  # Matches repeated words
```



1.4 Advanced Concepts

1.4.3 Common Regex Functions in Python

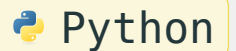
```
1  import re
2
3  text = "The price is $19.99"
4
5  # Different matching functions
6  re.search(r'\$\d+\.\d+', text)  # Finds first match
7  re.findall(r'\$\d+\.\d+', text) # Finds all matches
8  re.sub(r'\$(\d+\.\d+)', r'\1', text) # Substitution
9
10 # Splitting text
11 re.split(r'\s+', text) # Split on whitespace
```



1.5 Python Implementation

1.5.1 Basic Pattern Matching

```
1  import re
2  # Simple pattern matching
3  text = "The quick brown fox jumps over the lazy dog"
4  pattern = r"fox"
5  # Search for pattern
6  match = re.search(pattern, text)
7  if match:
8      print(f"Found '{pattern}' at position: {match.start()}-{match.end()}")
9  # Find all occurrences
10 words = re.findall(r"\w+", text)
11 print(f"All words: {words}")
```



1.5 Python Implementation

1.5.2 Email Validation Example


```
1  def is_valid_email(email):  
2      pattern = r'^[a-zA-Z0-9._%+-]+@[a-zA-Z0-9.-]+\.[a-zA-Z]{2,}$'  
3      return bool(re.match(pattern, email))  
4  # Test cases  
5  emails = [  
6      "user@example.com", # ✓  
7      "invalid.email@com", # ✗  
8      "user.name@bizerte.r-iset.tn", # ✓  
9      "@invalid.com" # ✗  
10 ]  
11 for email in emails:  
12     print(f"{email} → {'Valid' if is_valid_email(email) else 'Invalid'}")
```

 Python

1.5 Python Implementation

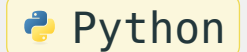
1.5.3 Phone Number Formatting

```
1  def format_phone_number(phone):
2      # Remove all non-digit characters
3      digits = re.sub(r'\D', '', phone)
4
5      # Format as (XXX) XXX-XXXX
6      if len(digits) == 10:
7          pattern = r'(\d{3})(\d{3})(\d{4})'
8          formatted = re.sub(pattern, r'(\1) \2-\3', digits)
9          return formatted
10
11     return "Invalid phone number"
```

 Python

1.5 Python Implementation

```
1  # Test cases
2  numbers = [
3      "1234567890",
4      "123-456-7890",
5      "(123) 456-7890",
6      "12345"
7  ]
8
9  for number in numbers:
10     print(f"{number} → {format_phone_number(number)}")
```




1.6 Best Practices

1. Use Raw Strings


- Always prefix regex patterns with `r` to avoid escape character issues

```
1 pattern = r'\d+' # Better than '\d+'
```

 Python

2. Compile Frequently Used Patterns

```
1 email_pattern = re.compile(r'^[\w\.-]+@[\w\.-]+\.\w+$')  
2 # Use multiple times  
3 email_pattern.match(email1)  
4 email_pattern.match(email2)
```

 Python


3. Be Specific

- Make patterns as specific as possible to avoid false matches;
- Use start (^) and end (\$) anchors when matching whole strings.

4. Test Thoroughly

- Test with both valid and invalid inputs
- Include edge cases in your tests


```
1 def test_pattern(pattern, test_cases):  
2     regex = re.compile(pattern)  
3     for test, expected in test_cases:  
4         result = bool(regex.match(test))  
5         print(f"'{test}': {'✓' if result == expected else 'x'}")
```

 Python

1.7 Common Pitfalls

1. Greedy vs. Non-Greedy Matching

```
1 # Greedy (default)
2 re.findall(r'<.*>', '<tag>text</tag>') # ['<tag>text</tag>']
3
4 # Non-greedy: Add (lazy) `?`
5 re.findall(r'<.*?>', '<tag>text</tag>') # ['<tag>', '</tag>']
```

 Python


2. Performance Considerations

- Avoid excessive backtracking (*recursion*);
- Be careful with nested quantifiers;
- Use more specific patterns when possible.

1.8 Applications


1. Basic Pattern Matching

```
1 # Write a pattern to match dates in format DD/MM/YYYY
2 date_pattern = r'\d{2}/\d{2}/\d{4}'
```

 Python

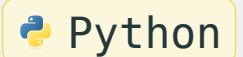
2. Data Extraction

```
1 # Extract all email addresses from text
2 text = "Contact us at support@example.com or sales@example.com"
3 emails = re.findall(r'[\w\.-]+@[\w\.-]+\.\w+', text)
```

 Python

3. Password Validation

```
1 def is_strong_password(password):  
2     # At least 8 chars, 1 upper, 1 lower, 1 digit, 1 special  
3     pattern = r'^(?=.*[A-Z])(?=.*[a-z])(?=.*\d)(?=.*[@$!%*?&])[A-Za-  
4         z\d@$!%*?&]{8,}$' # Positive Lookahead  
5     return bool(re.match(pattern, password))
```



1.8 Applications

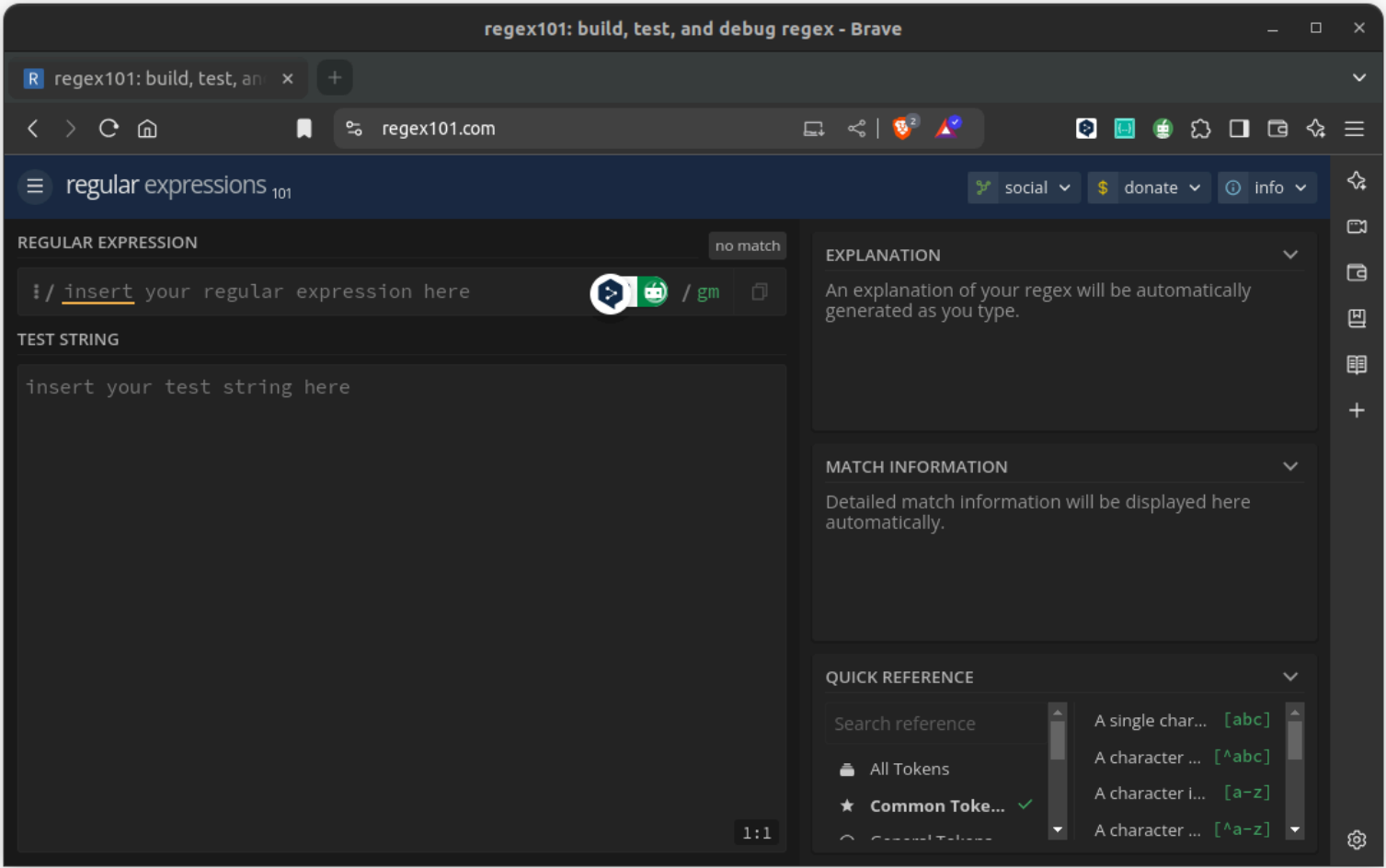


Figure 1: Build, test and debug regex patterns.

Thank you for your attention!