

✓ Introduction to Regular Expressions (RegEx) in Python

 Open in Colab

Regular Expressions (RegEx) help us **search, extract, and manipulate** text efficiently. Python provides the `re` module to work with regex.

In this notebook, we'll cover **simple regex patterns with diverse examples**, perfect for beginners!


```
1 import re # Import the regular expressions module
```

✓ 1. Searching for a Simple Character

Pattern: `a`

This pattern matches the letter 'a' anywhere in the string.

```
1 text = 'apple banana cat'
2 match = re.search(r'a', text)
3 if match:
4     print(f'Found "a" at position: {match.start()}') # Output: Position of first 'a'
5 else:
6     print('No match found.')
```


 Found "a" at position: 0

✓ 2. Matching an Exact Word

Pattern: `cat`

This pattern finds the word 'cat' in a string.

```
1 text = 'The cat sat on the mat.'
2 match = re.search(r'cat', text)
3 if match:
4     print(f'Found "cat" at position: {match.start()}')
5 else:
6     print('No match found.')
```

 Found "cat" at position: 4

✓ 3. Matching Any Character

Pattern: `c.t`

The `.` (dot) matches any single character. This finds words like 'cat', 'cot', or 'cut'.

```
1 text = 'cat cot cut'
2 matches = re.findall(r'c.t', text)
3 print(f'Matching words: {matches}') # Output: ['cat', 'cot', 'cut']
```

 Matching words: ['cat', 'cot', 'cut']

✓ 4. Matching a Digit

Pattern: `\d`

The `\d` pattern matches any single digit (0-9).

```
1 text = 'My number is 5 and yours is 7.'
2 matches = re.findall(r'\d', text)
3 print(f'Found digits: {matches}') # Output: ['5', '7']
```

 Found digits: ['5', '7']

5. Matching a Word Character

Pattern: `\w`

The `\w` pattern matches any letter, number, or underscore (`_`).

```
1 text = 'Hello, Python_3!'
2 matches = re.findall(r'\w', text)
3 print(f'Word characters: {matches}') # Output: ['H', 'e', 'l', 'l', 'o', ..., '3']
```

Word characters: ['H', 'e', 'l', 'l', 'o', 'P', 'y', 't', 'h', 'o', 'n', '_', '3']

6. Matching Whitespace

Pattern: `\s`

The `\s` pattern matches any space, tab, or newline character.

```
1 text = 'Hello World\nNew Line'
2 matches = re.findall(r'\s', text)
3 print(f'Whitespace characters found: {len(matches)}') # Output: 3 (space, newline, space)
```

Whitespace characters found: 3

7. Matching Multiple Occurrences

Pattern: `o+`

The `+` means 'one or more' repetitions of the letter 'o'.

```
1 text = 'soooooon moon good'
2 matches = re.findall(r'o+', text)
3 print(f'Found sequences: {matches}') # Output: ['oooo', 'oo', 'oo']
```

Found sequences: ['oooo', 'oo', 'oo']

8. Matching the Start of a String

Pattern: `^Hello`

The `^` symbol ensures the match is **only at the beginning** of the string.

```
1 text = 'Hello, world!'
2 match = re.match(r'^Hello', text)
3 if match:
4     print('The string starts with "Hello"!')
5 else:
6     print('No match at the start.')
```

The string starts with "Hello"!

9. Matching the End of a String

Pattern: `world!$`

The `$` symbol ensures the match is **only at the end** of the string.

```
1 text = 'This is my world!'
2 match = re.search(r'world!$', text)
3 if match:
4     print('The string ends with "world!"')
5 else:
6     print('No match at the end.')
```


The string ends with "world!"

10. Finding All Words with 'ing'

Pattern: `\b\w+ing\b`

This pattern finds full words ending in 'ing'.

```
1 text = 'I am running, singing, and playing today.'
2 matches = re.findall(r'\b\w+ing\b', text)
3 print(f'Words ending with "ing": {matches}') # Output: ['running', 'singing', 'playing']
```

 Words ending with "ing": ['running', 'singing', 'playing']

✓ Practical Example

The script below uses regular expressions to automatically extract essential business details including:


- email subjects
- recipient names
- dates
- monetary amounts
- company names

from sample email text, organizing the results into a structured pandas DataFrame for easy analysis.

```
1 import re
2
3 email = """Subject: Meeting Request for Product Launch
4     Dear John,
5
6     We are excited to announce the launch of our new product next month. The project cost is estimated at $12,500,000.
7     We would like to schedule a meeting on March 20th at 3 PM to discuss partnership opportunities.
8
9     Best regards,
10    Jane Doe
11    MarketingPro Inc."""
```


Extracting the Subject of the Email

```
1 # Pattern: r"Subject:\s*(.*)"
2 # - Matches the word "Subject:" followed by any whitespace (`\s*`).
3 # - Captures the rest of the line (the actual subject) using `(.*)`,
4 #   which means "any character (.) zero or more times (*)".
5 # - This helps in extracting the subject line from the email content.
6
7 match = re.search(r"Subject:\s*(.*)", email)
8 subject = match.group(1) if match else "Not Found"
9
10 print(subject)
```

 Meeting Request for Product Launch

Extracting the Recipient Name or Team

```
1 # Pattern: r"Dear\s+([\w\s]+),"
2 # - Matches the word "Dear" followed by one or more whitespace characters (`\s+`).
3 # - Captures the recipient's name or team using `([\w\s]+)`,
4 #   which means "one or more word characters (`\w`) or spaces (`\s`)".
5 # - Ensures the match ends before the comma (`,`), marking the end of the salutation.
6
7 match = re.search(r"Dear\s+([\w\s]+),", email)
8 recipient = match.group(1) if match else "Not Found"
9
10 print(recipient)
```

 John

Extracting Any Mentioned Date

```
1 # Pattern: r"(January|February|March|April|May|June|July|August|September|October|November|December) \d{1,2}"
2 # - Looks for any full month name (`January` to `December`).
3 # - Followed by a space and a **1 or 2 digit number** (`\d{1,2}`), representing the day.
4 # - This extracts dates like "March 20", "April 10", or "May 5" from the email content.
5
6 match = re.search(r"(January|February|March|April|May|June|July|August|September|October|November|December) \d{1,2}", email)
7 date_mentioned = match.group(1) if match else "Not Found"
8 print(date_mentioned)
```

➞ March

Extracting Monetary Values (Currency Amounts)

```
1 # Pattern: r"\$\d{1,3}(?:,\d{3})*(?:\.\d{2})?"
2 # - Matches a dollar sign (`$`).
3 # - Captures **1 to 3 digits** (`\d{1,3}`) at the start of the amount.
4 # - Optionally matches **comma-separated thousands** using `(?:,\d{3})*` (e.g., "$1,000" or "$15,000").
5 # - Optionally matches **decimal values** `(?:\.\d{2})?`, allowing for cents like "$99.99".
6 # - This extracts monetary values such as "$15,000", "$95,000", or "$499.99" from the email.
7
8 match = re.search(r"\$\d{1,3}(?:,\d{3})*(?:\.\d{2})?", email)
9 amount_mentioned = match.group(0) if match else "Not Found"
10 print(amount_mentioned)
```

➞ \$12,500,000

Extracting Company Names

```
1 # Pattern: r"\n\s*([\w\s]+(?:Inc\.|Ltd\.|Corp\.|Giant))"
2 # - Ensures the match starts on a new line (`\n`).
3 # - Allows for any leading whitespace (`\s*`).
4 # - Captures the company name using `([\w\s]+)`, which allows letters and spaces.
5 # - Ensures the company name ends with a **common business suffix** (e.g., "Inc.", "Ltd.", "Corp.", or "Giant").
6 # - This extracts names like "MarketingPro Inc.", "TechInnovators Ltd.", and "GlobalCorp."
7
8 match = re.search(r"^\\n\\s*([\\w \\s]+(?:Inc\\.|Ltd\\.|Corp\\.|Giant))", email, re.MULTILINE)
9 company_name = match.group(0) if match else "Not Found"
10 print(company_name)
11
```

➞ MarketingPro Inc.

Note:

If we want to match start-of-line within multi-line text we need to set the re.MULTILINE flag, otherwise it will match the first line only

The re.MULTILINE flag redefines ^ and \$ so that:

Anchor	Normal (default)	With re.MULTILINE
^	Start of string	Start of any line
\$	End of string	End of any line