

Text Processing For NLP

Understanding Regex

In this presentation, we will explore the power of regular expressions in Natural Language Processing (NLP) and how they can be applied to clean, analyze, and preprocess data.

Getting Started with Regex

Learn the basics of Regular Expressions, including patterns, characters, metacharacters, escaping, and character classes.



Patterns and Characters

Discover how to use regex to find specific patterns of characters in textual data.



Metacharacters and Escaping

Master the art of escaping special characters in regex and when to use metacharacters.



Character Classes and Ranges

Unleash the

power of character classes and ranges for matching different types of characters and numbers.

Patterns and Characters

- **Introduction to Patterns:** Regular expressions (regex) use patterns to manipulate text. These patterns consist of characters and metacharacters for specific search criteria.
- **Literal Characters:** Literal characters in a regex pattern match the exact characters themselves. For example, the pattern "apple" matches the word "apple".
- **Character Sets:** Character sets allow matching any one of a set of characters. For example, the pattern "[aeiou]" matches any vowel in the text.
- **Escaping Special Characters:** Special characters like ".", "\$", and "^" have special meanings in regex. To match them literally, they need to be escaped with a backslash (\).
- **Case Sensitivity:** By default, regex is case-sensitive. To perform case-insensitive matching, flags can be used in the regex pattern.

Metacharacters and Escaping

- **Metacharacters:** Metacharacters in regex have special meanings and for instance, "." matches any character, while "^" and "\$" respectively denote the start and end of a line.
- **Escaping Metacharacters:** To match metacharacters as literal characters, they must be escaped with a backslash. For example, to match the dot character ".", you would use "\."
- **Literal Matching:** By escaping metacharacters, you can precisely match specific patterns that would otherwise have special meanings.
- **Backslash Escaping:** Since the backslash itself is an escape character in many programming languages, it must be escaped as well when using regex. Use "\\\" to represent a single backslash.
- **Metacharacter Combinations:** Combining metacharacters with literal characters allows for powerful patterns for complex text manipulation.

Character Classes and Ranges

- **Character Classes:** Character classes, enclosed in square brackets "[]", match any one character from a set. For example, "[aeiou]" matches any one vowel.
- **Negation:** Using the caret symbol "^" within a character class negates the match. "[^0-9]" matches any non-digit character.
- **Character Ranges:** Ranges within character classes specify a range of characters. "[a-z]" matches any lowercase letter.
- **Combining Character Classes:** Character classes can be combined to create complex patterns. "[A-Za-z]" matches any letter, regardless of case.
- **Shorthand Character Classes:** Shortcuts like "\d" for digits, "\w" for word characters, and "\s" for whitespace simplify pattern creation.

Quantifiers and

Repetition

Quantifiers determine the number of times a preceding character should appear in the text. For instance, "a{3}" matches "aaa".

- **Asterisk (*)**: The asterisk quantifier matches zero or more occurrences of the preceding character or group. "ba*" matches "b", "ba", "baa", and so on.
- **Plus (+)**: The plus quantifier matches one or more occurrences of the preceding character or group. "ca+" matches "ca", "caa", "caaa", and so on.
- **Question Mark (?)**: The question mark quantifier matches zero or one occurrence of the preceding character or group. "da?" matches "d" and "da".
- **Braces ({ })**: Braces with a specific quantity, like "e{2}", match exactly the specified number of occurrences. "bee" matches "bee", but not "be".

Advanced Regex

Techniques

Explore advanced regular expression techniques, including anchors, boundaries, alternation, backreference, subpatterns, lookahead, and lookbehind.

Anchors and

Boundaries

Learn how to use anchors and

boundaries to match specific positions within a string of text.

Groups and

Capturing

Discover how to use groups

capturing to extract useful information.

Alternation and Logical

OR

Understand how to use alternation and

logical OR to match multiple patterns within a string.

Backreference and

Subpatterns

Master the art of backreference

subpatterns to match complex nested structures.

Groups and Capturing

- **Grouping with Parentheses:** Parentheses create groups to apply quantification to specific sections.
- **Capturing Groups:** Parentheses also capture the matched content, which can be accessed for extraction.
- **Named Capturing Groups:** Named groups provide a more descriptive reference to captured content.
- **Reusing Captured Content:** Captured content can be used later in the same pattern with backreferences.

Alternation and Logical

O.R. Alternation with |: The pipe symbol "|" allows multiple alternatives to be pattern.

- **Matching Multiple Alternatives:** For instance, "apple|banana" matches "banana".
- **Non-Capturing Groups (?:):** Parentheses with "?" after the opening parenthesis for non-capturing groups.
- **Balancing Options:** Alternation provides flexibility in matching different parts within a pattern.

Backreference and Subpatterns

- **Using Backreferences:** "\n" (where n is a number) matches content preceded by a group.
- **Repeating Captured Content:** Backreferences allow patterns like "(apple)" to match "applepieapple".
- **Nested Subpatterns:** Parentheses can be nested to create subpatterns for complex matches.
- **Subpattern Scope:** Subpatterns are useful for applying quantifiers and anchors to specific portions of the pattern.

Best Practices for Using Regex in Python

Learn how to use the re module in Python 3 to apply regex on text data and how to parse and extract information from real-world use cases.

Data Extraction with

Regex
Explore real-world use cases of regex in Python, including extracting URLs, emails, and phone numbers from unstructured text.

1

Introduction to re

Module
Get an overview of the re module in Python 3 and how to use it to apply regex on text data.

2

Cleaning and Preprocessing

Regex

Discover how to apply regex to preprocess and clean text prior to analysis.

3

Introduction to re

Module

Getting the Module: Begin by importing the re module in Python using

- **Using re.search():** Use `re.search()` to find the first match of a pattern in a string.
- **Using re.findall():** Employ `re.findall()` to extract all occurrences of a pattern in a string.
- **Flags for Flexibility:** Utilize flags like `re.IGNORECASE` for case-insensitive

Data Extraction with

Regexing URLs: Use regex to identify and extract URLs from text, aiding in data collection and analysis.

- **Capturing Emails:** Employ regex to capture and extract email addresses from documents.
- **Phone Number Extraction:** Regex assists in parsing and retrieving phone numbers from various formats.
- **Pattern Customization:** Adapt patterns to different data formats for accurate extraction.

Cleaning and Preprocessing with

Regex

Removing Unwanted Characters: Use regex to eliminate special characters, punctuation, or symbols.

- **Whitespace Management:** Replace multiple spaces with a single space for consistent formatting.
- **Text Normalization:** Apply regex for converting text to lowercase, standardizing representations.
- **Handling Redundancy:** Identify repeated characters or words with regex and replace them with a single instance.

Limitations and Best Practices of

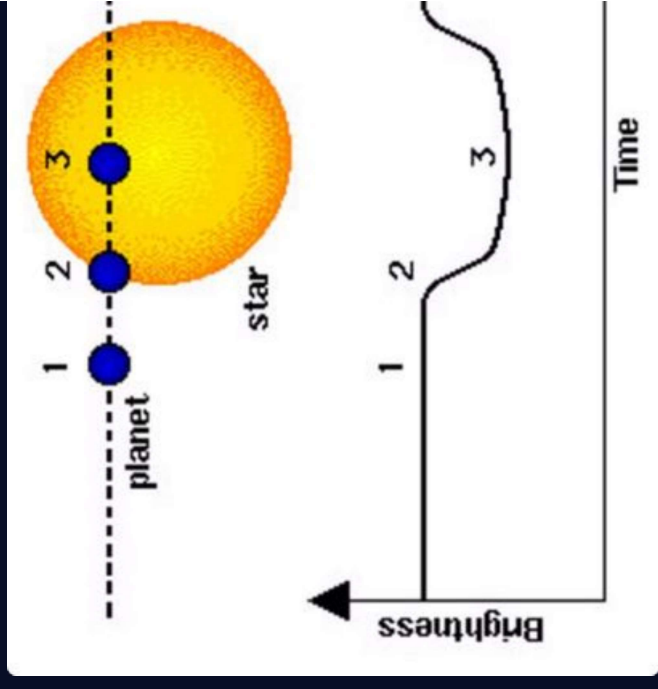
Understand the limitations of Regex and how to apply best practices to maximize performance.



Regex

Limitations

Explore the limits of regex when applied to Natural Language Processing and how to



Regex Best

Practices

Discover the best practices for applying regex to maximize performance and

Conclusion

Regular Expressions are essential for text processing and Natural Language Processing. Mastering the knowledge, skills, and best practices covered in this presentation, you will be able to use regex effectively and efficiently to your data processing needs.

Regex Basics

Patterns, Characters, Metacharacters, Escaping, Character Classes, Ranges, Quantifiers, Repetition.

Advanced Regex Techniques

Anchors, Boundaries, Groups, Capturing, Alternation, Logical OR, Backreference, Subpatterns, Lookahead, Lookbehind.

Regex I

Python

Cleaning a

Preproces

Best Pract