

Introduction to Regex and Web Scrapping

01 November 2023 08:32 PM

Regex or regex expression or in layman terms often called as a pattern is used when you want to extract a common pattern of data from a very large heap of organized/unorganized data.

Eg Data $\rightarrow \{1, \text{doc 1}, 2, \text{doc 2}, \text{doc 3}, \text{doc 4}, 2, 3, 5, 7, \dots\}$

Now the task is to find all correct document names with the correct number associated with them

\rightarrow If we observe carefully the common pattern for a correct document name: \rightarrow "doc" \rightarrow <number>

So the common pattern for the document name will be \rightarrow doc <number>

In order to apply regex, the input must be a string.

\Rightarrow Because of the above property, regex is widely used in NLP (A prime component of ML and DL)

At the very core regex is basically helping you verify the structure of a string by matching it with a pattern.

Real time use case of Regex

① Checking if input follows a pattern, like is the email correct

\hookrightarrow abc@gmail.com \checkmark
 \hookrightarrow abc#in \times

② In real time data it often happens that 2 different words which mean the same thing need to be captured.

Eg
iphone \swarrow Apple phone
Both refer to the phone only

③ Extracting specific portion of a text. Eg. a text column containing the address is given, and we want to extract postal code.

Mr John Smith. 132, My Street, Kingston, New York 12401

Name

Street No.

Area

City

Pincode

The above break down can be done using Regex.

④ Retrieve/replace/structure

\hookrightarrow find particular word/structure or replace a word detected by pattern.

⑤ Splitting a text based on the occurrence of a specific character.

Creating Regex Patterns in Python

Package \rightarrow import re

\rightarrow Understanding patterns

\rightarrow Literal

\hookrightarrow When you give exactly what you want in your pattern

\hookrightarrow Text \rightarrow maabnabcab
Pattern \rightarrow ab

\hookrightarrow It matches literally, i.e. case sensitive

\rightarrow Metachar

\hookrightarrow special characters with a general representative meaning

\hookrightarrow These characters can be combined together to create complex patterns to extract specific text.

Depending upon the requirement literal & metachar can also be combined together to create a pattern

eg var|wt