

# Pattern Matching using Regular Expression

## Pattern matching :

Pattern matching, involves the use of one-dimensional string matching.

Patterns are either tree structures or sequences.

There are different classes of programming languages and machines which make use of pattern matching.

In the case of machines, the major classifications include deterministic finite state automata, deterministic pushdown automata, nondeterministic pushdown automata and Turing machines.

Regular programming languages make use of regular expressions for pattern matching.

Tree patterns are also used in certain programming languages like Haskell as a tool to process data based on the structure.

Compared to regular expressions, tree patterns lack simplicity and efficiency.

There are many applications for pattern matching in computer science.

High-level language compilers make use of pattern matching in order to parse source files to determine if they are syntactically correct.

In programming languages and applications, pattern matching is used in identifying the matching pattern or substituting the matching pattern with another token sequence.

## Regular Expression :

A regular expression in a programming language is a special text string used for describing a search pattern.

It is extremely useful for extracting information from text such as code, files, log, spreadsheets or even documents.

While using the regular expression the first thing to recognize is that everything is essentially a character, and we are writing patterns to match a specific sequence of characters also referred as string.

Ascii or latin letters are those that are on your keyboards and Unicode is used to match the foreign text.

It includes digits and punctuation and all special characters like \$#@!%, etc.

Regular expression could tell a program to search for specific text from the string and then to print out the result accordingly. Expression can include

- Text matching
- Repetition
- Branching
- Pattern-composition etc.

In Python, a regular expression is denoted as RE (REs, regexes or regex pattern) are imported through re module.

Python supports regular expression through libraries.

In Python regular expression supports various things like Modifiers, Identifiers, and White space characters.

To learn regular expressions we can use below sites

<https://regexr.com/>

The screenshot displays the regexr.com website. The browser's address bar shows the URL <https://regexr.com>. The website's interface is dark-themed. On the left, there is a 'Menu' sidebar with links to 'Pattern Settings', 'My Patterns', 'Cheatsheet', 'RegEx Reference', 'Community Patterns', and 'Help'. The main content area is divided into sections: 'Expression' (showing the regex `/([A-Z])\w+/g`), 'Text' (containing a paragraph about regexr.com), and 'Tools' (which includes a 'Replace' button and a detailed explanation of the regex components). The 'Tools' section explains that the expression consists of a 'Capturing group #1' (the parentheses around `A-Z`), a 'Character set' (the square brackets around `A-Z`), and a 'Word' (the `\w` character class).

<https://regex101.com/>

The screenshot displays the regex101.com website. The browser's address bar shows the URL <https://regex101.com>. The website's interface is light-themed. On the left, there is a sidebar with sections: 'SAVE & SHARE' (with a 'Save RegEx' button), 'FLAVOR' (with radio buttons for PCRE (PHP), ECMAScript (JavaScript), Python, and Golang), and 'TOOLS' (with links to 'Code Generator' and 'Regex Debugger'). The main content area is divided into sections: 'REGULAR EXPRESSION' (with an input field), 'TEST STRING' (with an input field), and 'SUBSTITUTION' (with an output field). On the right, there is an 'EXPLANATION' section that provides a detailed breakdown of the regex components, including 'MATCH INFORMATION' and a 'QUICK REFERENCE' table.

## Important concepts of regular expressions

### Metacharacters :

Metacharacters are characters with a special meaning

Character	Description	Example
[]	A set of characters	"[a-m]"
\	Signals a special sequence (can also be used to escape special characters)	"\d"
.	Any character (except newline character)	"he..o"
^	Starts with	"^hello"
\$	Ends with	"world\$"
*	Zero or more occurrences	"aix*"
+	One or more occurrences	"aix+"
{}	Exactly the specified number of occurrences	"al{2}"
	Either or	"falls stays"
()	Capture and group	

### Special Sequences :

A special sequence is a \ followed by one of the characters in the list below, and has a special meaning:

Character	Description	Example
\A	Returns a match if the specified characters are at the beginning of the string	"\AThe"
\b	Returns a match where the specified characters are at the beginning or at the end of a word	r"\bain" r"ain\b"
\B	Returns a match where the specified characters are present, but NOT at the beginning (or at the end) of a word	r"\Bain" r"ain\B"
\d	Returns a match where the string contains digits (numbers from 0-9)	"\d"
\D	Returns a match where the string DOES NOT contain digits	"\D"
\s	Returns a match where the string contains a white space character	"\s"
\S	Returns a match where the string DOES NOT contain a white space character	"\S"
\w	Returns a match where the string contains any word characters (characters from a to Z, digits from 0-9, and the underscore _ character)	"\w"
\W	Returns a match where the string DOES NOT contain any word characters	"\W"
\Z	Returns a match if the specified characters are at the end of the string	"Spain\Z"

## Sets :

A set is a set of characters inside a pair of square brackets [] with a special meaning

Set	Description
[arn]	Returns a match where one of the specified characters ( <b>a</b> , <b>r</b> , or <b>n</b> ) are present
[a-n]	Returns a match for any lower case character, alphabetically between <b>a</b> and <b>n</b>
[^arn]	Returns a match for any character EXCEPT <b>a</b> , <b>r</b> , and <b>n</b>
[0123]	Returns a match where any of the specified digits ( <b>0</b> , <b>1</b> , <b>2</b> , or <b>3</b> ) are present
[0-9]	Returns a match for any digit between <b>0</b> and <b>9</b>
[0-5][0-9]	Returns a match for any two-digit numbers from <b>00</b> and <b>59</b>
[a-zA-Z]	Returns a match for any character alphabetically between <b>a</b> and <b>z</b> , lower case OR upper case
[+]	In sets, <b>+</b> , <b>*</b> , <b>.</b> , <b> </b> , <b>()</b> , <b>\$</b> , <b>{}</b> has no special meaning, so <b>[+]</b> means: return a match for any <b>+</b> character in the string

