

Agenda

- Regular Expressions
 - Simple pattern matching
 - `.` character
 - `\` character
 - Metacharacters
 - Matching digits (and non digits) - `\d` , `\D`
 - Matching word (and non-word) characters - `\w` , `\W`
 - Matching whitespace characters - `\s` , `\S`
 - Anchors
 - Word boundary (and non-boundary) anchors - `\b` , `\B`
 - Beginning and end anchors - `^` , `$`
 - Character Set
 - Range notation - `[a-z]` , `[A-Z]` , `[1-6]`
 - Negation - `[^a-z]`
 - Quantifiers
 - `*` , `+` , `?` , `{3}` , `{3, 4}`
 - Groups
 - `()` and `|`
 - Functions in `re` library
-

Regex

Regular Expression

➤ shorthand notations for:

- ⇒ Searching
- ➔ parsing
- ➔ manipulating

Complex Patterns in String

Q: why regex?

- Help in NLP
- Extracting Information from Unstructured Text Data

Business Use-Cases

- Extracting customer info from feedbacks
- Validating Texts such as Email
- Masking PII data such as Phone-numbers, address etc.

Ex: Email Validation

* Note

- ① Regex is Case sensitive
- ② Regex is Order sensitive

Meta Characters

Special characters that don't match themselves

- ① $\backslash d$: Matches all digits
- ② $\backslash D$: Matches everything except digits
- ③ $\backslash w$: Matches Alphanumeric and $_$ (Under Score)
- ④ $\backslash W$: Matches Non-Alphanumeric
- ⑤ $\backslash s$: Matches white space characters (space, tab, newline)
- ⑥ $\backslash S$: Matches non white space characters

* Note : Uppercase counterparts of meta characters 'Negate' the lower case ones

Other Special Characters

.	- Any Character Except New Line
\d	- Digit (0-9)
\D	- Not a Digit (0-9)
\w	- Word Character (a-z, A-Z, 0-9, _)
\W	- Not a Word Character
\s	- Whitespace (space, tab, newline)
\S	- Not Whitespace (space, tab, newline)

} meta characters

\b	- Word Boundary
\B	- Not a Word Boundary
^	- Beginning of a String
\$	- End of a String

Anchors

[]	- Matches Characters in brackets
[^]	- Matches Characters NOT in brackets
	- Either Or
()	- Group

Documentation: <https://docs.python.org/3/howto/regex.html>

Anchors

- * Anchors don't match with any character
- * They match with positions instead

① \b : Word boundary

↳ new-line or space is considered as boundary

② ^ : Matches only the beginning of a string

③ \$: Matches only the end of string

Character Set

Matches for one of the characters defined inside set

Ex

[ab]

[.-]

* numeric character set

[123456]

[1-6]

[1-9]



* Alphabet character

[a-z]

[A-Z]

[a-zA-Z]

* Negation in character-set

[^a-z] ⇒ match everything but lowercase char

Quantifiers and Groups

Quantifiers: matching more than one char at a time

① $?$: 0 or One match

② $+$: 1 or more

③ $*$: 0 or more

④ $\{n\}$ \Rightarrow Exactly n times

$\{a, b\}$ Range of num (min, max)
 $\{3, 4\}$

Groups

$(x|s|rs)$

re-functions

- ① **match** : Checks for a match for input at the beginning of string
- ② **search** : Finds only the first Occ. of a given Expression
- ③ **findall** : Find all the Occ. of a given Expression
- ④ **finditer** : Finds all occ. one by one as an iterable

* For string manipulation

- ① **sub** : Used for substitution operation
Ex: `re.sub('ab', 'ba', 'abab')`
output: babac
- ② **split** : Used for splitting the input string into list

re Flags

- ① `re.IGNORECASE` or `re.I`
 - ① Ignore case sensitivity

- ② `re.VERBOSE` or `re.V`
 - ① It allows for spacing indentation and comments while writing Expression
 - ② You can use this when you have highly non-understandable complex pattern.