



Data Science Cheat Sheet

Python Regular Expressions

SPECIAL CHARACTERS

- ^** | Matches the expression to its right at the start of a string. It matches every such instance before each **\n** in the string.
- \$** | Matches the expression to its left at the end of a string. It matches every such instance before each **\n** in the string.
- .** | Matches any character except line terminators like **\n**.
- ** | Escapes special characters or denotes character classes.
- A|B** | Matches expression **A** or **B**. If **A** is matched first, **B** is left untried.
- +** | Greedily matches the expression to its left 1 or more times.
- *** | Greedily matches the expression to its left 0 or more times.
- ?** | Greedily matches the expression to its left 0 or 1 times. But if **?** is added to qualifiers (**+**, *****, and **?** itself) it will perform matches in a non-greedy manner.
- {m}** | Matches the expression to its left **m** times, and not less.
- {m,n}** | Matches the expression to its left **m** to **n** times, and not less.
- {m,n}?** | Matches the expression to its left **m** times, and ignores **n**. See **?** above.

CHARACTER CLASSES

[A.K.A. SPECIAL SEQUENCES]

- \w** | Matches alphanumeric characters, which means **a-z**, **A-Z**, and **0-9**. It also matches the underscore, **_**.
- \d** | Matches digits, which means **0-9**.
- \D** | Matches any non-digits.
- \s** | Matches whitespace characters, which include the **\t**, **\n**, **\r**, and space characters.
- \S** | Matches non-whitespace characters.
- \b** | Matches the boundary (or empty string) at the start and end of a word, that is, between **\w** and **\W**.
- \B** | Matches where **\b** does not, that is, the boundary of **\w** characters.

\A | Matches the expression to its right at the absolute start of a string whether in single or multi-line mode.

\Z | Matches the expression to its left at the absolute end of a string whether in single or multi-line mode.

SETS

- []** | Contains a set of characters to match.
- [amk]** | Matches either **a**, **m**, or **k**. It does not match **amk**.
- [a-z]** | Matches any alphabet from **a** to **z**.
- [a\ -z]** | Matches **a**, **-**, or **z**. It matches **-** because **** escapes it.
- [a-]** | Matches **a** or **-**, because **-** is not being used to indicate a series of characters.
- [-a]** | As above, matches **a** or **-**.
- [a-z0-9]** | Matches characters from **a** to **z** and also from **0** to **9**.
- [+*?)]** | Special characters become literal inside a set, so this matches **(**, **+**, *****, and **)**.
- ^ab5]** | Adding **^** excludes any character in the set. Here, it matches characters that are not **a**, **b**, or **5**.

GROUPS

- ()** | Matches the expression inside the parentheses and groups it.
- (?)** | Inside parentheses like this, **?** acts as an extension notation. Its meaning depends on the character immediately to its right.
- (?PAB)** | Matches the expression **AB**, and it can be accessed with the group name.
- (?aiLmsux)** | Here, **a**, **i**, **L**, **m**, **s**, **u**, and **x** are flags:
 - a** — Matches ASCII only
 - i** — Ignore case
 - L** — Locale dependent
 - m** — Multi-line
 - s** — Matches all
 - u** — Matches unicode
 - x** — Verbose

(?:A) | Matches the expression as represented by **A**, but unlike **(?PAB)**, it cannot be retrieved afterwards.

(?#...) | A comment. Contents are for us to read, not for matching.

A(?:B) | Lookahead assertion. This matches the expression **A** only if it is followed by **B**.

A(?:!B) | Negative lookahead assertion. This matches the expression **A** only if it is not followed by **B**.

(?<=B)A | Positive lookbehind assertion.

This matches the expression **A** only if **B** is immediately to its left. This can only match fixed length expressions.

(?<!B)A | Negative lookbehind assertion.

This matches the expression **A** only if **B** is not immediately to its left. This can only match fixed length expressions.

(?P=name) | Matches the expression matched by an earlier group named "name".

(...)\1 | The number **1** corresponds to the first group to be matched. If we want to match more instances of the same expression, simply use its number instead of writing out the whole expression again. We can use from **1** up to **99** such groups and their corresponding numbers.

POPULAR PYTHON RE MODULE FUNCTIONS

- re.findall(A, B)** | Matches all instances of an expression **A** in a string **B** and returns them in a list.
- re.search(A, B)** | Matches the first instance of an expression **A** in a string **B**, and returns it as a re match object.
- re.split(A, B)** | Split a string **B** into a list using the delimiter **A**.
- re.sub(A, B, C)** | Replace **A** with **B** in the string **C**.

re.match object | `match = re.search(A, B)`
`match[0]` = entire match
`match[1]` = first capture group
`match[2]` = second capture group

re.compile | `pat = re.compile('pattern')`
`pat.search(B)`, `pat.split(B)`, etc