Mission Name : Advanced Regular Expression.

* <u>Capture groups</u> :- we define a capture group by wrapping the part we want to capture in parentheses-

Eg:-
```
pattern = r"(\w+SQL)"
sql-flavors = titles.str.extract(pattern, flags= re.I)
sql-flavours_freq = sql-flavors.value_counts()
print(sql-flavors_freq)
```

* Another example of capture groups:-

Eg:-
```
pattern = r"[Pp]ython ([\d\.]+)"
py-versions = titles.str.extract(pattern)
py-versions_freq = dict(py-version.value_counts())
```

The above code will extract the version followed by the word Python/python and make a frequency dictionary with it.

* Eg:- pattern = r"\b[Cc][^\.\+\w]\b"

* Lookarounds:-

| Lookaround | Pattern | Explanation |
|---|---|---|
| Positive Lookahead | zzz (?=abc) | Matches only when zzz is followed by abc. |
| Negative Lookahead | zzz (?!abc) | Matches only when zzz is not followed by abc. |
| Positive lookbehind | (?<=abc)zzz | Match zzz only if it is preceded by abc |
| Negative lookbehind | (?<!abc)zzz | Match zzz only if it is NOT preceded by abc. |

Eg:- pattern = r"(?<!series\s)\b[Cc]\b(?![\.\+])"

## * Back References:-

If we want to repeat any capture group in a RegEx pattern, we can use back references.

We can number capture groups from left to right and just mention the number wherever we want to repeat it.

Eg:- (Hello)(Goodbye)\2\1

This will match HelloGoodbyeGoodbyeHello.

(\w)\1

with match 2 repeating characters, like "ee", "oo" etc.

Eg 2:- pattern = r"\b(\w+)\s\1\b"

## * re.sub()

Function is similar to string.replace().

Syntax:- re.sub(pattern, repl, string, flags=0)

regex   string to   string
        replace

For pandas series:- Series.str.replace(pat, repl, flags=0)

Eg:- titles_clean = titles.str.replace(pat = r"e-?\s?mail",
                                        repl = "email",
                                        flags = re.I)

* Extracting domain names:- pattern = r".{4}s?://([\w\.\-]+)/?"
* Naming extracted columns:- r"(?P<date>.+) (?P<time>.+)"