

همان طور که گفته شد از ریجکس می توان برای جستجو، جایگزینی و تطابق الگو استفاده کرد. در این راستا به معرفی سه تابع این کتابخانه می پردازیم:

تابع search :

```
txt="some text"
x=re.search("your pattern",txt)
```

خروجی قطعه کد بالا شامل یک رشته و یک تاپل است که به ترتیب عبارت تطابق داده شده و محل شروع و پایان عبارت مذکور را در صورت وجود بر می گرداند. اگر جستجو نتیجه ای نداشته باشد و عبارت مورد نظر در متن نباشد خروجی None خواهد بود.

حال فرض کنیم می خواهیم وجود یک عبارت یا یک الگو را در رشته دلخواه بررسی کنیم. برای مثال فرض کنیم می خواهیم ببینیم آیا کلمه "Python" در یک رشته وجود دارد یا خیر:

```
import re

string = "Python is fun"

match = re.search('Python', string)

if match:
    print("pattern found inside the string")
else:
    print("pattern not found")

# Output: pattern found inside the string
```

تابع find all :

```
txt="some text"
x=re.findall("your pattern",txt)
```

دستور بالا یک لیست شامل تمامی عبارت هایی که در متن با الگو تطابق دارند را خروجی می دهد.

برای مثال فرض کنیم می خواهیم همه اعداد استفاده شده در یک رشته را

Regex در پایتون



عطیه سادات طباطبایی
دانشجو علوم کامپیوتر
دانشگاه تهران
atiehtaba82@gmail.com

Regex، مخفف عبارت Regular Expression، دنباله ای از کاراکترها است که یک الگوی جستجو را تعریف می کنند. ریجکس یک ابزار قدرتمند برای پردازش متن و تطبیق الگو است. زبان های برنامه نویسی مانند java، Python و Perl از ریجکس پشتیبانی می کنند. این کتابخانه معمولاً در اعتبارسنجی داده ها، تجزیه متن و عملیات جستجو و جایگزینی استفاده می شود.

ریجکس در واقع ترکیبی از کاراکترهای خاص است که الگویی را تعریف می کنند که با مجموعه خاصی از کاراکترها مطابقت دارد.

در ادامه ابتدا به معرفی قواعد الگوسازی ریجکس می پردازیم و سپس چند مثال را با هم می بینیم.

در جدول زیر برخی از پرکاربردترین قواعد این کتابخانه را مشاهده می کنید:

General Tokens	
.	Any character
\n	Newline character
\t	Tab character
\s	Any whitespace character (Including \t, \n, etc)
\S	Any non-whitespace character
\w	Any word character (Upper/lowercase letters, 0-9, _)
\W	Any non-word character
\b	Word boundary (Matches between characters)
\B	Non-word boundary
^	The start of a line
\$	The end of a line
\\	The literal character "\"

Quantifiers	
a b	Match either "a" or "b"
?	Match either "a" or "b"
+	One or more
*	Zero or more
*?	Zero or more, but stop after first match
{N}	Exactly N number of times (Where N is number)
{N, M}	From N to M number of times (Where N and M are numbers)

برای مثال "[A-Za-z]+" یک الگو است که با هر رشته ای که فقط از حروف تشکیل شده مطابقت دارد.

برای استفاده از ریجکس در زبان پایتون از دستور زیر استفاده می کنیم.

```
import re
```

که نشان دهنده سال است و دو جفت دوعدد که به ترتیب نشان دهنده ماه و روز هستند و اعداد هر قسمت با کاراکتر "/" از هم جدا شده اند مانند ۲۰۲۳/۰۵/۱۰. می خواهیم اگر کاربر تاریخ تولدش را مطابق الگوی خواسته شده وارد نکند ارور دریافت کند. خواهیم داشت :

```
import re
user_birthdate=input()
x=re.search("[0-9]{4}/[0-9]{2}/[0-9]{2}",user_birthdate)

if not x:
    print("pleas enter your birthdate correctly")
```

اعتبار سنجی آدرس ایمیل:

فرض کنید از کاربر خواسته ایم تا آدرس ایمیل خود را وارد کند. می خواهیم ببینیم آیا عبارتی که کاربر وارد کرده می تواند یک آدرس ایمیل باشد؟

برای این کار ابتدا یک الگو تعریف می کنیم. اگر الگوی تعریف شده با عبارت وارد شده توسط کاربر مطابقت داشت رشته ورودی می تواند یک آدرس ایمیل باشد در غیر این صورت عبارت وارد شده معتبر نخواهد بود.

برای تعریف الگو ابتدا به ساختار یک آدرس ایمیل دقت می کنیم. در یک آدرس ایمیل ابتدا مجموعه ای از کاراکترها که می توانند شامل عدد، حروف و یا کاراکترهای دیگر باشند کنار هم قرار می گیرند سپس علامت "@" می آید و پس از آن نیز مجددا مجموعه ای از کاراکترها می آید و در آخر نیز دامنه سایت که با "." شروع می شود :

```
import re
address=input()
x=re.search("[A-Za-z0-9._-]+@[A-Za-z0-9._-]+\.[A-Za-z0-9._-]+",address)
if x:
    print("valid email address")
else:
    print("invalid email address")
```

دقت کنید چون "." جزء قواعد ریجکس است و برای تطابق هر کاراکتر دلخواه به شمار می رود، اگر بخواهیم از آن به عنوان خود کاراکتر "." استفاده کنیم قبل از آن از علامت "\" استفاده می کنیم.

سخن پایانی :

کار با ریجکس در نگاه اول می تواند دشوار به نظر برسد اما با تمرین می توانید از قابلیت های این کتابخانه قدرتمند بهره ببرید.

برای مطالعه بیشتر می توانید به آدرس https://www.w3schools.com/python/python_regex.asp رجوع کنید.

استخراج کنیم؛ با توجه جدول بالا الگوی مورد نظر ما "[۰-۹]+" خواهد بود بنابراین طبق دستور زیر عمل می کنیم:

```
import re

string = 'The first release candidate of Python 3.12 was offered on 6 August 2023.'
pattern = '[0-9]+'

result = re.findall(pattern, string)
print(result)

# Output: ['3', '12', '6', 2023]
```

تابع sub :

```
txt="some text"
x=re.sub("your pattern","new string",txt)
```

با استفاده از دستور بالا می توانید یک عبارت جدید را جایگزین الگوی تعریف شده کنید.

برای مثال فرض کنیم می خواهیم همه فاصله های یک رشته را حذف کنیم. برای این کار باید همه فاصله ها را با یک رشته خالی جایگزین کنیم:

```
import re

string = 'Hello World!'
# matches all whitespace characters
pattern = '\s+'
# empty string
replace = ''

new_string = re.sub(pattern, replace, string)
print(new_string)

# Output: HelloWorld!
```

حال برای درک بهتر کاربرد regex چند مثال را با هم بررسی می کنیم

مطابقت فرمت تاریخ (yyyy/mm/dd):

فرض کنیم از کاربر خواسته ایم تا تاریخ تولدش را مطابق فرمت بالا وارد کند. فرض می کنیم فقط فرمت بالا را به عنوان تاریخ می پذیریم (چهار عدد