

## Pattern matching and Web Scraping

It is used to describe a search pattern,

when we want to extract any pattern from the given raw data, for eg: mobile No., email id, client id etc.

```
In [2]: 1 #Searching a pattern in string 'findall'
2 #findall returns list of matches
3 import re
4 Nameage = 'David is 25 and Smith is 30 /n Michael is 28 and Wayne is 35'
5 ages = re.findall('\d{1,2}',Nameage)
6 print (ages)
7 names = re.findall('[A-Z][a-z]*',Nameage)
8 print (names)
9 agedict ={}
10 x=0
11 for i in names:
12     agedict[i]=ages[x]
13     x+=1
14 print (agedict)
```

```
['25', '30', '28', '35']
['David', 'Smith', 'Michael', 'Wayne']
{'David': '25', 'Smith': '30', 'Michael': '28', 'Wayne': '35'}
```

```
In [4]: 1 #we can directly search for a string in a given string using 'search'
2 #search returns a single match
3 if re.search('match','I want to match the string'):
4     print ('match captured')
```

match captured

```
In [9]: 1 # #to get the index range of the pattern match
2 str_ = 'This is a demo regex prog for a regex understanding'
3 for i in re.finditer('regex',str_):
4     print (i)
5     index = i.span()
6     print (index)
7
```

```
<re.Match object; span=(15, 20), match='regex'>
(15, 20)
<re.Match object; span=(32, 37), match='regex'>
(32, 37)
```

```
In [15]: 1 #compile method which catches patterns and provide method to substitute
2 demo = 'Java html c++ ruby html'
3 object_ = re.compile('html') #matching objects with compile
4 sub_ = object_.sub('python',demo)
5 sub_
```

Out[15]: 'Java python c++ ruby python'

```
In [17]: 1 num = '123 1234 12345 123456 1234567 87654321'
2 print ('Matches :', len(re.findall(r'\d{5,7}' , num))) #use len to give the
Matches : 4
```

## Web Scrapping

Scrap useful data from web and store it in csv format or excel format.

```
In [33]: 1 #Zomato customer care India
2 import urllib.request #importing urllib package to read data from url
3 import re
4 url = 'http://www.talkingtrends.com/zomato-customer-care-number-address-cont
5 response=urllib.request.urlopen(url) #opening the url
6 html = response.read() #reading the url
7 htmlstr=html.decode() #decoding the url
8 data=re.findall('\d{3} - \d{8}',htmlstr) #Matching the pattern in the above
9 for i in data:
10     print (i)
```

```
079 - 60601010
080 - 60601010
044 - 60601010
011 - 60601010
040 - 60601010
030 - 60601010
022 - 60601010
020 - 60601010
141 - 60601010
079 - 60601010
080 - 60601010
044 - 60601010
011 - 60601010
040 - 60601010
030 - 60601010
022 - 60601010
020 - 60601010
141 - 60601010
```

**practise questions:**

```
In [5]: 1 # from string x = 'my name is Michael and my no. is +919865471232' extract '  
2 # search for 'john' in the given string  
3 # modify the phone no. with some other no.  
4 #  
5 # import re  
6 # x = 'my name is Michael and my no. is +919865471232'  
7 # name = re.findall('[A-Z][a-z]*',x)  
8 # phone = re.findall ('\+91\d{10}',x)  
9  
10 # nameNo={}  
11 # j=0  
12 # for i in name:  
13 #     nameNo[i]=phone[j]  
14 # print (nameNo)  
15  
16 # if re.search('John',x):  
17 #     print ('match captured')  
18 # else:  
19 #     print ('No match')  
20  
21 obj = re.compile('\+91\d{10}')  
22 obj1 = obj.sub('kjhkjf',x)  
23 print (obj1)  
24  
25 # for i in re.finditer('Michael',x):  
26 #     index = i.span()  
27  
28 # print (i)
```

my name is Michael and my no. is kjhkjf

```
In [77]: 1 s = 'Rahul is learning python'  
2 name = re.compile('[A-Z][a-z]*')  
3 new_s = name.sub('Michael',s)  
4 new_s  
5
```

Out[77]: 'Michael is learning python'

```
In [ ]: 1 import re
```

```
In [ ]: 1 David is 25 and Smith is 30 /n Michael is 28 and Wayne is 35
```