

Regular Expression

```

1 import re
2 string = "The quick brown fox jumps over the lazy dog."
3 pattern="brown"
4 match=re.search(pattern, string) # returns True/False, used for validation
5 if match:
6     print("Match found:", match.group())
7 else:
8     print("Match not found")

```

➞ Match found: brown

```

1 import re
2 string = "The quick brown fox jumps over the lazy dog."
3 pattern="The"
4 match=re.match(pattern, string)
5 if match:
6     print("Match found:", match.group())
7 else:
8     print("Match not found")

```

Match found: The

```

1 import re
2 string = "The quick brown fox jumps over the lazy dog."
3 pattern="[aeiou]"
4 matches=re.findall(pattern, string)
5 if matches:
6     print("Match found:", matches)
7 else:
8     print("Match not found")

```

Match found: ['e', 'u', 'i', 'o', 'o', 'u', 'o', 'e', 'e', 'a', 'o']

```

1 import re
2 string = "The quick brown fox jumps over the lazy dog."
3 pattern="brown"
4 new_str=re.sub(pattern, "red", string)
5 print(new_str)
6 print(string) # actual string

```

The quick red fox jumps over the lazy dog.
The quick brown fox jumps over the lazy dog.

```

1 import re
2 string = "The quick brown fox jumps over the lazy dog."
3 pattern="\s+"
4 split_str=re.split(pattern, string)
5 print(split_str)

```

['The', 'quick', 'brown', 'fox', 'jumps', 'over', 'the', 'lazy', 'dog.']

Match any character: . (dot) Example: a.b matches "a b", "a1b", "a!b", etc.

Match the start of a string: ^ Example: ^hello matches "hello world", but not "world hello"

Match the end of a string: *Example* : *world* matches "hello world", but not "world hello"

Match any single digit: \d Example: \d matches any digit from 0 to 9

Match any non-digit character: \D Example: \D matches any character that is not a digit

Match any whitespace character: \s Example: \s matches spaces, tabs, and newlines

Match any non-whitespace character: \S Example: \S matches any character that is not whitespace

Match any word character (letters, digits, or underscore): \w Example: \w+ matches one or more word characters (letters, digits, or underscore)

Match any non-word character: \W Example: \W matches any character that is not a word character

Match zero or one occurrence of the previous character: ? Example: colou?r matches "color" and "colour"

Match zero or more occurrences of the previous character: * Example: ab*c matches "ac", "abc", "abbc", etc.

Match one or more occurrences of the previous character: + Example: ab+c matches "abc", "abbc", "abbbc", etc.

Match a specific number of occurrences of the previous character: {m} Example: a{3} matches "aaa"

Match a range of occurrences of the previous character: {m,n} Example: a{2,4} matches "aa", "aaa", and "aaaa"

Match any one of a set of characters: [] Example: [abc] matches "a", "b", or "c"

Match any character except those in a set: [^] Example: [^abc] matches any character except "a", "b", or "c"

Match the same characters as a previously captured group: \1, \2, etc. Example: (\w) is \1 matches "a is a", "b is b", etc.

These are just some of the most common regular expressions used in Python. There are many more advanced regular expressions available as well.

▼ validation through Regular Expression (RE)

1. re.search(pattern, string) ---- to check if pattern is present or not
2. re.findall(pattern, string) ---- to find all occurrences of pattern in string
3. re.sub(pattern, newpattern, inp) to replace string
- 4.

```
1 import re
2 p='[a-z]'
```

3 string=input("Enter data to check: ")

```
4 if re.search(p, string):
5     print("yes")
6 else:
7     print("No")
8
9 p2='\d'
```

10 string=input("Enter data to check: ")

```
11 if re.search(p2, string):
12     print("yes")
13 else:
14     print("No")
```

```
Enter data to check: amar123
yes
Enter data to check: amar123panchal
yes
```

```
1 import re
2 p='[a-zA-Z0-9]'
```

```
3 string=input("Enter data to check: ")
4 if re.search(p, string):
5     print("yes")
6 else:
7     print("No")
```

```
Enter data to check: Amar123panchal
yes
```

```
1 import re
2 p='(.com|.co.in|.in|man)'
```

```
3 i=input("Enter data to check: ")
4 if re.search(p, i):
5     print("yes")
6 else:
7     print("No")
```

```
Enter data to check: google.com
yes
```

```
1 import re
2 p='[a-zA-Z0-9]+\.' # seaches for dot at end
```

```
3 string=input("Enter data to check: ")
4 if re.search(p, string):
5     print("yes")
6 else:
7     print("No")
```

```
Enter data to check: asds121cs.
yes
```

```
1 import re
2 p='[0-9]+\.[0-9]+\.[0-9]+\.[0-9]' # IP validation
```

```
3 string=input("Enter data to check: ")
4 if re.search(p, string):
5     print("yes")
6 else:
7     print("No")
```

```
Enter data to check: 123amar.hehe.gone
No
```

```
1 import re
2 p='[a-zA-Z0-9]+\@[a-zA-Z]+\.[a-zA-Z]' # email validation
3 string=input("Enter data to check: ")
4 if re.search(p, string):
5     print("yes")
6 else:
7     print("No")
```

Enter data to check: sachin@iyla.com
yes

```
1 import re
2 p='\d{10}' # mobile validation
3 i=input("Enter data to check: ")
4 if re.search(p, i):
5     print("yes")
6 else:
7     print("No")
```

Enter data to check: 8787878787
yes

```
1 import re
2 p='^d{2}' # starts with 2
3 i=input("Enter data to check: ")
4 if re.search(p, i):
5     print("yes")
6 else:
7     print("No")
```

Enter data to check: 2656
yes

```
1 import re
2 p='d{2}$' # ends with 2
3 i=input("Enter data to check: ")
4 if re.search(p, i):
5     print("yes")
6 else:
7     print("No")
```

Enter data to check: 5662
yes

```
1 import re
2 p='d{2}$' #
3 i=input("Enter data to check: ")
4 if re.search(p, i):
5     print("yes")
6 else:
7     print("No")
```

▼ Extraction

```

1 pattern="[\\W]+" #word separated
2 inp=input("Enter data to check: ")
3 data=re.findall(pattern, inp)
4 print("Final result: ", data)

```

```

Enter data to check: ajkda 789 jnk** k*(
Final result: [' ', ' ', '** ', '(']

```

```

1 pattern="\+(\d{2})-(\d{4})-(\d{6})" # pattern for telephone number
2 newpattern=r"Country:\1 Area:\2 Number\3" # r"" is raw string, will replace the matching pattern
3 inp=input("Enter data to check: ")
4 newstr=re.sub(pattern, newpattern, inp)
5 print("Final result: ", newstr)

```

```

Enter data to check: +91-1923-220007
Final result: Country:91 Area:1923 Number220007

```

Double-click (or enter) to edit

```

1 import re
2 log=''192.168.198.92 - - [22/Dec/2002:23:08:37 -0400] "GET
3   / HTTP/1.1" 200 6394 www.yahoo.com
4   "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1...)" "-"
5 192.168.198.92 - - [22/Dec/2002:23:08:38 -0400] "GET
6   /images/logo.gif HTTP/1.1" 200 807 www.yahoo.com
7   "http://www.some.com/" "Mozilla/4.0 (compatible; MSIE 6...)" "-"
8 192.168.72.177 - - [22/Dec/2002:23:32:14 -0400] "GET
9   /news/sports.html HTTP/1.1" 200 3500 www.yahoo.com
10  "http://www.some.com/" "Mozilla/4.0 (compatible; MSIE ...)" "-"
11 192.168.72.177 - - [22/Dec/2002:23:32:14 -0400] "GET
12  /favicon.ico HTTP/1.1" 404 1997 www.yahoo.com
13  "-" "Mozilla/5.0 (Windows; U; Windows NT 5.1; rv:1.7.3)..." "-"
14 192.168.72.177 - - [22/Dec/2002:23:32:15 -0400] "GET
15  /style.css HTTP/1.1" 200 4138 www.yahoo.com
16  "http://www.yahoo.com/index.html" "Mozilla/5.0 (Windows..." "-"
17 192.168.72.177 - - [22/Dec/2002:23:32:16 -0400] "GET
18  /js/ads.js HTTP/1.1" 200 10229 www.yahoo.com
19  "http://www.search.com/index.html" "Mozilla/5.0 (Windows..." "-"
20 192.168.72.177 - - [22/Dec/2002:23:32:19 -0400] "GET
21  /search.php HTTP/1.1" 400 1997 www.yahoo.com
22  "-" "Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; ...)" "-"
23

```

```

1 pattern="[0-9]+\.[0-9]+\.[0-9]+\.[0-9]"
2 iplist=re.findall(pattern, log)
3 print(iplist)

```

```

['192.168.198.9', '192.168.198.9', '192.168.72.1', '192.168.72.1', '192.168.72.1', '192.168.72.1', '192.168.72.1']

```

```

1 pattern="[0-9]{2}/[a-zA-Z]{3}/[0-9]{4}"
2 datelist=re.findall(pattern, log)

```

```
3 print(datelist)

['22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002']

1 pattern="\d+/\w{3}/\d+"
2 datelist=re.findall(pattern, log)
3 print(datelist)

['22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002', '22/Dec/2002']
```

▼ Web-Scraping

```
1 from bs4 import BeautifulSoup
2 import requests
3 f=requests.get("https://www.timesjobs.com/candidate/job-search.html?searchType=personalizedSearch&from=submit&txtKeywords=python&txtLocation=mumbai&cboworkExp1=0").text
4
5 soup=BeautifulSoup(f, 'lxml')
6 print(soup)

<!DOCTYPE html>
<html><head>
<link href="https://fonts.googleapis.com/css?family=Poppins:400,500,600,700" rel="stylesheet"/>
<link href="https://fonts.googleapis.com/icon?family=Material+Icons" rel="stylesheet"/>
<link href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/materialize.css?v=7.1.7" media="all" rel="stylesheet" type="text/css"/>
<link href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/global-usability.css?v=7.1.7" media="all" rel="stylesheet" type="text/css"/>
<link href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/srp-usability.css?v=7.1.7" media="all" rel="stylesheet" type="text/css"/>
<script src="https://static.timesjobs.com/newtj_js/scripts/tj_scripts/usability/jquery-3.3.1.min.js" type="text/javascript"></script>
<script src="https://static.timesjobs.com/newtj_js/scripts/jquery.tokeninput.js" type="text/javascript"></script>
<!-- <script type="text/javascript" src="https://code.angularjs.org/1.7.0/angular.min.js"></script> -->
<script src="https://static.timesjobs.com/newtj_js/scripts/tj_scripts/usability/angular1.7.0.min.js" type="text/javascript"></script>
<script src="https://static.timesjobs.com/newtj_js/scripts/tj_scripts/usability/materialize.min.js" type="text/javascript"></script>
<script src="https://static.timesjobs.com/newtj_js/scripts/groupJs/dwrGroup.js?v=7.1.7" type="text/javascript"></script>
<script src="https://static.timesjobs.com/newtj_js/scripts/tj_scripts/mustache.js" type="text/javascript"></script>
<link href="https://www.timesjobs.com/candidate/job-search.html" rel="canonical"/>
<link href="https://www.timesjobs.com/candidate/job-search.html?searchType=personalizedSearch&from=submit&txtKeywords=python&txtLocation=mumbai&cboworkExp1=0&sequence=2&sequence=2" rel="canonical"/>
<title>python Jobs in mumbai with Entry Level experience - TimesJobs</title>
<meta content="Jobs in python in mumbai with Entry Level experience. Search and apply jobs online at Timesjobs.com" name="description"/>
<!-- Google Tag Manager -->
<script>(function(w,d,s,l,i){w[l]=w[l]||[];w[l].push({'gtm.start':
new Date().getTime(),event:'gtm.js'});var f=d.getElementsByTagName(s)[0],
j=d.createElement(s),dl1!=&#39;dataLayer&#39;?&#39;l&#39;+&#39;1&#39;:'';j.async=true;j.src=
'https://www.googletagmanager.com/gtm.js?id='+i+dl1&#39;f.parentNode.insertBefore(j,f);
})(window,document,'script','dataLayer','GTM-59SJDG4');</script>
<!-- End Google Tag Manager -->
<script src="https://www.gstatic.com/charts/loader.js" type="text/javascript"></script>
<script>
//google.charts.load('current', {'packages':['bar']});
google.charts.load('current', {'packages: ['corechart', 'bar']});
</script>
<script>!function(e){var n="https://s.go-mpulse.net/boomerang/";if("False"=="True")e.BOOMR_config=e.BOOMR_config||{},e.BOOMR_config.PageParams=e.BOOMR_config.PageParams||{},e.BOOMR_config.PageP
<body>
<div id="site">
<!-- Header Starts -->
<!-- Google Tag Manager (noscript) -->
<noscript>
<iframe height="0" src="https://www.googletagmanager.com/ns.html?id=GTM-59SJDG4" style="display:none;visibility:hidden" width="0"></iframe>
```

```

</noscript>
<!-- End Google Tag Manager (noscript) -->
<meta content="text/html; charset=utf-8" http-equiv="Content-Type"/>
<script src="https://wchat.freshchat.com/js/widget.js" type="text/javascript"></script>
<script>
window.fcWidget.init({
token: "d7130618-a917-400b-804c-bdf44d8d178b",
host: "https://wchat.freshchat.com",
config: {
headerProperty: {
hideChatButton: true
}
}
});

</script>
<!--[if IE]>
<script src="scripts/html5shiv.js"></script>
<![endif]-->
<script src="https://static.timesjobs.com/newtj_css/css/tj_css/usability/global-usability.css?v=7.1.7" type="text/css"></script>

```

```

1 from bs4 import BeautifulSoup
2 import requests
3 f=requests.get("https://www.timesjobs.com/candidate/job-search.html?searchType=personalizedSearch&from=submit&txtKeywords=python&txtLocation=mumbai&cboworkExp1=0").text
4
5 soup=BeautifulSoup(f, 'lxml')
6 soup.prettify()
7 #print(soup)

```

```

'<!DOCTYPE html>\n<html>\n <head>\n  <link href="https://fonts.googleapis.com/css?family=Poppins:400,500,600,700" rel="stylesheet"/>\n  <link href="https://fonts.googleapis.com/icon?family=Material+Icons" rel="stylesheet"/>\n  <link href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/materialize.css?v=7.1.7" media="all" rel="stylesheet" type="text/css"/>\n  <link href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/global-usability.css?v=7.1.7" media="all" rel="stylesheet" type="text/css"/>\n  <link

```

```

1 from bs4 import BeautifulSoup
2 import requests
3 f=requests.get("https://www.timesjobs.com/candidate/job-search.html?searchType=personalizedSearch&from=submit&txtKeywords=python&txtLocation=mumbai&cboworkExp1=0").text
4
5 soup=BeautifulSoup(f, 'lxml')
6 soup.prettify().split()
7
8 #print(soup)

```

```

['<!DOCTYPE',
'html>',
'<html>',
'<head>',
'<link',
'href="https://fonts.googleapis.com/css?family=Poppins:400,500,600,700"',
'rel="stylesheet"/>',
'<link',
'href="https://fonts.googleapis.com/icon?family=Material+Icons"',
'rel="stylesheet"/>',
'<link',
'href="https://static.timesjobs.com/newtj_css/css/tj_css/usability/materialize.css?v=7.1.7"',
'media="all"',

```

```
1 from bs4 import BeautifulSoup
2
3 html_doc = """
4 <html>
5 <head>
6 <title>My Website</title>
7 </head>
8 <body>
9 <div id="content">
10 <h1>Welcome to my website!</h1>
11 <p>This is a paragraph of text.</p>
12 <a href="http://www.google.com">Google</a>
13 </div>
14 </body>
```



```
15 </html>
16 """
17 soup=BeautifulSoup(html_doc, 'html.parser')
18 div_tag=soup.find('div', {'id':'content'})
19 print(div_tag)

<div id="content">
<h1>Welcome to my website!</h1>
<p>This is a paragraph of text.</p>
<a href="http://www.google.com">Google</a>
</div>
```

```
1 from bs4 import BeautifulSoup
2 html_doc = """
3 <html>
4 <head>
5 <title>My Website</title>
6 </head>
7 <body>
8 <div class="article">
9 <h2>Article Title</h2>
10 <p>This is the first paragraph.</p>
11 <p>This is the second paragraph.</p>
12 </div>
13 <div class="article">
14 <h2>Another Article Title</h2>
15 <p>This is another paragraph.</p>
16 </div>
17 </body>
18 </html>
19 """
20 soup=BeautifulSoup(html_doc, 'html.parser')
21 h2tags=soup.find_all('h2')
22 for h2tag in h2tags:
23     print(h2tag)

<h2>Article Title</h2>
<h2>Another Article Title</h2>
```

1

