



# NLP

## Lab-1

### Regular Expressions

#### Lab Objectives

- Learn about regular expressions
- Use python to implement regular expressions through a library named (re)

#### Regular expressions (Regex) (RE)

A regular expression is a sequence of characters that define a search pattern. Regular expression is used in search engines, search and replace dialogs of word processors and text editors, and in lexical analysis.

#### Syntax of regular expressions

Each character in a regular expression is either a metacharacter, having a special meaning, or a regular character that has a literal meaning.

#### Metacharacters in python re library

Meta character	Description	Examples
*	<b>zero or more:</b> Matches any number of occurrences of the previous characters	<b>ad*e matches</b> ae, ade, adde, addde,...
?	<b>zero or one:</b> Matches at most one occurrence of the previous characters	<b>ad?e matches</b> ae and ade
+	Matches <b>one or more</b> occurrences of the previous characters	<b>ad+e matches</b> ade, adde, addde, ...
+, *, ?	The '*', '+', and '?' qualifiers are all greedy; they match as much text as possible. Adding ? matches <b>as few characters as possible</b>	<b>Given &lt;a&gt; b &lt;c&gt;</b> <b>&lt;.*&gt; matches &lt;a&gt; b &lt;c&gt;</b> <b>&lt;.*?&gt; matches &lt;a&gt;</b>
{n}	Matches <b>exactly n</b> occurrences of the previous RE	<b>ad{2}e matches</b> adde
{n,}	<b>n or more:</b> Matches n or more occurrences of the previous RE	<b>ad{2,}e matches</b> Adde, addde, ...
{n,m}	<b>from n to m:</b> Matches from n to m occurrences of the previous characters	<b>ad{2,4}e matches</b> Adde, addde, adddde.
{n,m}?	<b>fewer characters as possible</b> will be matched	
.	Matches <b>one occurrence of any character</b> of the alphabet or digit except new line (can include new line character with the re.DOTALL flag)	<b>a.e matches</b> aae, aAe, abe, aBe, a1e, etc.

ad and e are constants  
(regular characters)



Meta character	Description	Examples
<b>.</b>	Matches <b>any string of characters</b> and until it encounters a new line character	
<b>[...]</b>	Matches <b>one occurrence of any character contained in the list</b>	<b>[abc]</b> means one occurrence of either a, b, or c
<b>[^...]</b>	Matches <b>one occurrence of any character not contained in the list</b>	<b>[^abc]</b> means one occurrence of any character that is not an a, b, and c
<b>[A-Z]</b>	Matches <b>one upper-case letter</b>	
<b>[a-z]</b>	Matches <b>one lower-case letter</b>	
<b>[0-9]</b>	Matches <b>one digit</b>	
<b>\d</b>	Matches <b>one digit</b> . Equivalent to [0-9]	<b>A\dC</b> matches A0C, A1C, A2C, A3C,...
<b>\D</b>	Matches <b>one non-digit</b> . Equivalent to [^0-9]	
<b>\w</b>	Match <b>one word character</b> : letter, digit, or underscore. Equivalent to [a-zA-Z0-9_]	<b>1\w2</b> matches 1a2, 1A2, 1b2, 1B2, ...
<b>\W</b>	Matches <b>one special character</b> . Equivalent to [^\w]: non alphanumeric @! # ...	
<b>\s</b>	Matches <b>one whitespace character</b> : space, tabulation, new line, form feed, etc.	
<b>\S</b>	Matches <b>one non-white space character</b> . Equivalent to [^\s]	
<b> </b>	Or operator between two RE	<b>a bc</b> matches a or bc
<b>()</b>	Grouping of a pattern	<b>(a b)c</b> matches ac or bc
<b>^</b>	Match <b>from the start of string</b> (or <b>start of line</b> with re.MULTILINE flag)	<b>^cat\$</b> matches cat alone in a line
<b>\$</b>	Match <b>the end of string or line</b>	
<b>\</b>	Match <b>special character</b>	Match * using \*
<b>\b</b>	<b>To match whole word</b> . Matches the empty string, but only at the beginning or end of a word	<b>'\bfoo\b'</b> matches 'foo', 'foo.', '{foo}', 'bar foo baz' but not 'foobar' or 'foo3'

## Python re library

- To import re library we use the statement: `import re`
- 3 main functions are provided:
  - `re.search()`: checks for a match to the pattern anywhere in the string
  - `re.match()`: checks for a match to the pattern only at the beginning of the string
  - `re.findall()`: returns list of all matches non-overlapping in the string



- Each function takes 3 parameters:
  - Pattern: the regex
  - String: the string that will be searched to find the pattern
  - Flag: The expression's behavior can be modified by specifying a flags value

## Examples

	Code	output
1	<pre>m= re.search('.*','abc133\\33') #Match any string print(m.group(0))</pre>	abc133\33
2	<pre>m= re.search('.*?','abc133\\33') #Note the change when adding ? print(m) print(m.group(0)) #prints nothing</pre>	<_sre.SRE_Match object; span=(0, 0), match="">>
3	<pre>m= re.search('abcd*','abc133\\33') #Match any string has abc and d with any number of occurrences print(m.group(0))</pre>	abc
4	<pre>m= re.search('abcd+','abc133\\33') #Match any string has abc and d. d should appear 1 or more print(m) print(m.group(0)) #Runtime error</pre>	None
5	<pre>m= re.search('abcd+',' abc31abcdd33\\33') #Match any string has abc and d. d should appear 1 or more print(m.group(0))</pre>	Abcdd
6	<pre>m= re.search('^abcd+',' abc31abcdd33\\33') #Match any string has abc and d appear in the beginning. #d should appear 1 or more print(m) m.group(0) #Runtime error</pre>	None
7	<pre>m= re.search('^a{3}',' abc31abcdd33\\33') #Match any string has aaa appear in the beginning. print (m)</pre>	None
8	<pre>m= re.search('^a{3,5}','aaaabc31abcdd33\\33') #Match any string has aaa or aaaa or aaaaa appear in the beginning. print (m)</pre>	<_sre.SRE_Match object; span=(0, 4), match='aaaa'>
9	<pre>m= re.search('a{3,5} b{2} ^c{2}','aaaabc31abcdd33\\33') #Match any string has aaa or aaaa or aaaaa or bb or cc in the beginning. print (m)</pre>	<_sre.SRE_Match object; span=(0, 4), match='aaaa'>



	Code	output
10	<pre>m= re.search('[fek]', 'aaaabc31fabcd33\\3') #Match f or e or k. print (m)</pre>	<_sre.SRE_Match object; span=(0, 4), match='f'>

## match.group(n)

Returns one or more subgroups of the match.

**Example:**

```
m = re.match(r"(\w+) (\w+)", "Isaac Newton, physicist")
m.group(0)      # The entire match
#'Isaac Newton'
m.group(1)      # The first parenthesized subgroup.
#'Isaac'
m.group(2)      # The second parenthesized subgroup.
#'Newton'
m.group(1, 2)   # Multiple arguments give us a tuple.
#('Isaac', 'Newton')
```

## Alternative method to define pattern

If the pattern will be used multiple times on different sequences, it is define a regex object storing the pattern and then use this object to call the match, search, or other methods.

**Example:**

<pre>import re reObj = re.compile('\d+')#find all numbers m= reObj.findall('aaaabc31fabcd33\\3') print (m)</pre>	<pre>['31', '33', '3']</pre>
--	------------------------------

## Tasks

- (1) Search on Vodafone number on the string where the length number should be 11 and start with 010. The String is 'My phone number is 01033192192'
- (2) Find All adverbs in this text 'He was carefully disguised but captured quickly by police'
- (3) Split The string based on number '20A50B1C19D', hint: try split function



## Solutions

(1) Vodafone number:

```
re.search('010[0-9]{8}', 'My phone number is 01033192192.')  
print (m.group(0))
```

(2) Adverbs

```
text = "He was carefully disguised but captured quickly by  
police."  
m = re.findall(r"\w+ly", text)  
print (m)
```

(3) Split by numbers

```
m = re.split('\d+', '20A34B13C')  
print (m)
```

## Student task

Write a regex to extract all emails from a text.

## For more Information

<https://docs.python.org/2/library/re.html>