

# Regular Expressions in Python

L435/L555  
Dept. of Linguistics, Indiana University  
Fall 2014

## Regular Expression Module

### Module

In order to use regular expressions, we need to load the module.

```
import re
```

## Regular Expression Symbols

.	wildcard
\	escapes specials characters
[...]	character set
[^...]	complement of character set
	or
*	Kleene star: 0 or more (of previous)
+	Kleene plus: 1 or more (of previous)
{ m,n }	repeat between m and n times
^	beginning of a string
\$	end of a string

## Understanding regular expressions

See slides 28–37 here: <http://cl.indiana.edu/~md7/13/245/slides/04-searching/slides.pdf>

## Regular Expression Functions

<code>compile(&lt;pattern&gt;)</code>	compiles a regex pattern into a pattern object – for reuse
<code>search(&lt;pattern&gt;, &lt;string&gt;)</code>	searches for regex pattern in string
<code>match(&lt;pattern&gt;, &lt;string&gt;)</code>	checks at <b>beginning</b> of string
<code>split(&lt;pattern&gt;, &lt;string&gt;)</code>	splits the string based on pattern, returns a <b>list</b>
<code>findall(&lt;pattern&gt;, &lt;string&gt;)</code>	returns a list of all occurrences
<code>sub(&lt;pat&gt;, &lt;rep&gt;, &lt;string&gt;)</code>	replaces pat by rep in string

## Example

```
import re
```

```
mysent = input('Gimme a sentence!\n')  
if (not re.search('[\u!\\.?:]', mysent)):  
    print('this is not a sentence')
```

Example

```
import re

mysent = input('Gimme a sentence!\n')
newstr = re.sub('[A-Z]', 'XX', mysent)
print(newstr)
```

Pattern Objects

Module

The functions `compile`, `search`, and `match` return a pattern object. The objects contain information about the pattern itself and for the matching functions also information about the matched segments in the string.

```
import re

phoneNums = re.compile('^(?\\d{3}[-\\s])\\d{3}[-\\s]\\d{4}$')
myphone = input('Give me a phone number: ')
if phoneNums.search(myphone):
    print('format correct')
else:
    print('format incorrect')
```

Example

```
import re

mysent = 'a rose is a rose is a rose'
allstr = re.search('(\\s)', mysent)
print(allstr.group(1))

allstr = re.findall('(\\s)', mysent)
print(allstr)
```