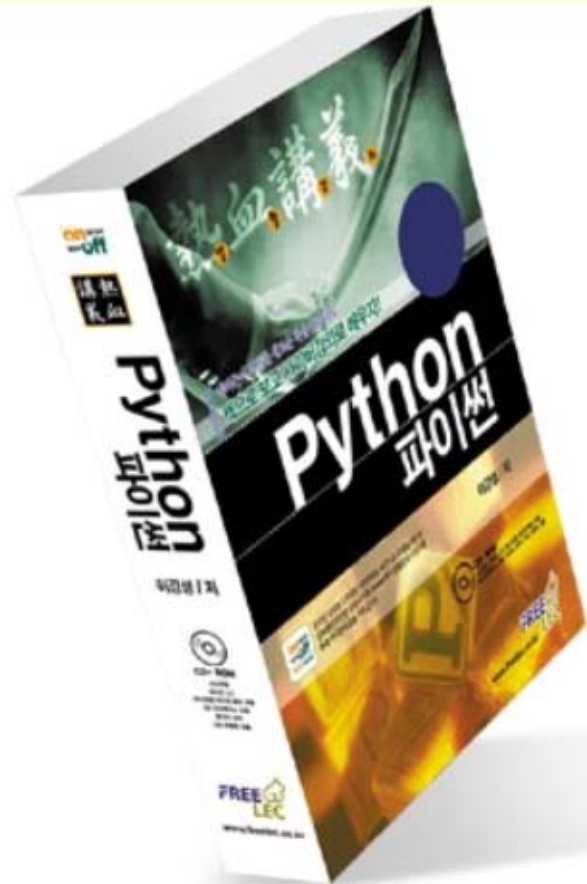


# 熱血講義

프리렉의 열혈강의 시리즈

## Python 파이썬



# Python

❖ 20

:

( gslee@mail.gwu.ac.kr )  
2



1.

2.

3.

4.

5.

6.

re



/

**‘ 876time a123 3 1.234 .567 892. a.b123’**



**(Regular Expression) ?**



**, ,**



**- re**

## (meta) ?



*	0	ca*t    ct, cat, caat, caaaaaat
+	1	ca+t    cat, caaaaaat    .
?	0          1	ca?t    ct, cat    .
{m}	m	ca{2}    caa    .
{m,n}	m          n	ca{2,4}    caat, caaat, caaaat    .

.	re.DOTALL
^	<ol style="list-style-type: none"> <li>1.</li> <li>2. re.MULTILINE</li> <li>3. [] . [^5] 5가</li> </ol>
\$	<ol style="list-style-type: none"> <li>1.</li> <li>2. [] \$</li> </ol>
[]	<p>[abc] "a", "b", "c" . [a-c]</p> <p>- . [a-zA-Z0-9]</p>
	a   b a (or) b
()	

# match

## - match

➤ `match(pattern, string[, flags])`

➤ `:` **Match**  
**None**

```
>>> re.match('[0-9]', '1234')  
<SRE_Match object at 00B06640>  
>>> re.match('[0-9]', 'abc') #  
>>>
```

# match

## ● More examples

```
>>> m = re.match('[0-9]', '1234') #
```

```
>>> m.group()
```

```
'1'
```

```
>>> m = re.match('[0-9]+', '1234') # 1
```

```
>>> m.group()
```

```
'1234'
```

```
>>> m = re.match('[0-9]+', '1234 ') #
```

```
>>> m.group()
```

```
'1234'
```

```
>>> re.match('[0-9]+', '    1234 ') # .
```

```
>>> re.match(r'[\t]*[0-9]+', '    123') # 가
```

```
<_sre.SRE_Match object at 0x00AEEA60>
```



\\	.
\\d	. [0-9]
\\D	가 . [^0-9]
\\s	. [ \\t\\n\\r\\f\\v]
\\S	가 . [^ \\t\\n\\r\\f\\v]
\\w	. [a-zA-Z0-9_]
\\W	가 . [^a-zA-Z0-9_]
\\b	. 가 .
\\B	\\b 가 .

```
>>> m = re.match('\\s*\\d+', '    123')
>>> m.group()
'    123'
>>> m = re.match('\\s*(\\d+)', '    1234 ')
>>> m.group(1) #
'1234'
```

# search

## – search

- match –

- search –



- `search(pattern, string[, flags])`

```
>>> re.search('\d+', ' 1034 ').group()
```

```
'1034'
```

```
>>> re.search('\d+', ' asf123 ').group()
```

```
'123'
```

## ● raw

```
>>> '\d' #
```

```
'\\d'
```

```
>>> '\s' #
```

```
'\\s'
```

```
>>> '\b' #
```

```
'\x08'
```

```
>>> '\\b' #          \      가
```

```
'\\b'
```

```
>>> r'\b'          # raw
```

```
'\\b'
```

## ● raw

```
>>> '\\\' # \
```

```
'\\'
```

```
>>> r'\\\' # \
```

```
'\\\\'
```

```
# \ \
```

```
>>> re.search('\\\\\\', ' \ ' ).group()
```

```
'\\'
```

```
>>> re.search(r'\\', ' \ ' ).group()
```

```
'\\'
```

## ● raw

```
>>> re.search('\d+', ' !123! ').group()
'123'
>>> re.search('\d+', ' abc123 ').group()
'123'
>>> re.search('\b\d+\b', ' !123! ')
>>> re.search('\b\d+\b', ' abc123 ')
>>> re.search('\\b\\d+\\b', ' !123! ').group()
'123'
>>> re.search('\\b\\d+\\b', ' abc123 ')
```

- (greedy)
- 가
- (non-greedy)

*?	* .
+?	+ .
??	? .
{m,n}?	{m,n} .



```
>>> s = '<a href="index.html">HERE</a><font size="10">'
>>> re.search(r'href="(.*)"', s).group(1) # greedy
'index.html">HERE</a><font size="10'

>>> re.search(r'href="(.*?)"', s).group(1)
'index.html'
```

## ➤ match

<code>group()</code>	.
<code>start()</code>	.
<code>end()</code>	.
<code>span()</code>	( , ) .

`match`



## ➤ match group

<code>group( )</code>	
<code>group(n)</code>	n . 0
<code>groups( )</code>	

`group`

## ● group

```
>>> m = re.match(r'(\d+)(\w+)(\d+)',  
'123abc456')
```

```
>>> m.group()
```

```
'123abc456'
```

```
>>> m.group(0)
```

```
'123abc456'
```

```
>>> m.group(1)
```

```
'123'
```

```
>>> m.group(2)
```

```
'abc45'
```

```
>>> m.group(3)
```

```
'6'
```

```
>>> m.groups()
```

```
('123', 'abc45', '6')
```

```
>>> m.start()
```

```
0
```

```
>>> m.end()
```

```
9
```

```
>>> m.span()
```

```
(0, 9)
```

## ● group

```
>>> m = re.match('(a(b(c))(d))',  
'abcd')
```

```
>>> m.group(0)
```

```
'abcd'
```

```
>>> m.group(1)
```

```
'abcd'
```

```
>>> m.group(2)
```

```
'bc'
```

```
>>> m.group(3)
```

```
'c'
```

```
>>> m.group(4)
```

```
'd'
```

- `(?P<name>...)`
- `name : , ... :`

```
>>> m = re.match(r'(?P<var>[_a-zA-Z]\w*)\s*=\s*(?P<num>\d+)', 'a = 123')
```

```
>>> m.group('var')
```

```
'a'
```

```
>>> m.group('num')
```

```
'123'
```

```
>>> m.group()
```

```
'a = 123'
```

<code>compile(pattern[, flags])</code>	pattern ( ).
<code>search(pattern, string[, flags])</code>	string pattern 가 .
<code>match(pattern, string[, flags])</code>	string pattern .
<code>split(pattern, string[, maxsplit = 0])</code>	string pattern .
<code>findall(pattern, string)</code>	string pattern .
<code>sub(pattern, repl, string[, count = 0])</code>	string pattern repl . (20-5 )
<code>subn(pattern, repl, string[, count = 0])</code>	sub . (20-5 )
<code>escape(string)</code>	가 .

## ● `re.split(pattern, string[, maxsplit = 0])`

```
>>> re.split('\W+', 'applt, orange and spam.') #
      ,      ,      1
```

```
['applt', 'orange', 'and', 'spam', '']
```

```
>>> re.split('(\W+)', 'applt, orange and spam.') #      가
      .
```

```
['applt', ',', ' ', 'orange', ' ', ' ', 'and', ' ', ' ', 'spam', '.', '']
```

```
>>> re.split('\W+', 'applt, orange and spam.', 1) # maxsplit
      가
      .
```

```
['applt', 'orange and spam.']
```

```
>>> re.split('\W+', 'applt, orange and spam.', 2)
```

```
['applt', 'orange', 'and spam.']
```

```
>>> re.split('\W+', 'applt, orange and spam.', 3)
```

```
['applt', 'orange', 'and', 'spam.']
```

## ● `findall(pattern, string)`

```
>>> re.findall(r'[_a-zA-Z]\w*', '123 abc 123 def')
['abc', 'def']

>>> s = '''
<a href="link1.html">link1</a>
<a href="link2.html">link2</a>
<a href="link3.html">link3</a>'''

>>> re.findall(r'href="(.*?)"', s)
['link1.html', 'link2.html', 'link3.html']
```

- **compile(pattern[, flags])**



가

```
>>> p = re.compile(r'([_a-zA-Z]\w*)\s*=\s*(\d+)')
```

```
>>> m = p.match('a = 123')
```

```
>>> m.groups()
```

```
('a', '123')
```

```
>>> m.group(1)
```

```
'a'
```

```
>>> m.group(2)
```

```
'123'
```



## ● compile 가

<b>I, IGNORECASE</b>	.
<b>L, LOCALE</b>	\w, \W, \b, \B
<b>M, MULTILINE</b>	^가 \$ ^, \$
<b>S, DOTALL</b>	.(\\n)
<b>U, UNICODE</b>	\w, \W, \b, \B가
<b>X, VERBOSE</b>	#

```
>>> import re
>>> p = re.compile('the', re.I)
>>> p.findall('The cat was hungry. They were scared
because of the cat')
['The', 'The', 'the']
```

<code>search(string[, pos[, endpos]])</code>	string . pos, endpos string .
<code>match(string[, pos[, endpos]])</code>	string .
<code>split(string[, maxsplit = 0])</code>	re.split .
<code>findall(string)</code>	re.findall .
<code>sub(repl, string[, count = 0])</code>	re.sub .
<code>subn(repl, string[, count = 0])</code>	re.subn .

- `search(string[, pos[, endpos]])`

```
>>> t = 'My ham, my egg, my spam'
```

```
>>> p = re.compile('(my)', re.I)
```

```
>>> pos = 0
```

```
>>> while 1:
```

```
    m = p.search(t, pos)
```

```
    if m:
```

```
        print m.group(1)
```

```
        pos = m.end()
```

```
    else: break
```

```
My
```

```
my
```

```
my
```

- **sub(repl, string[, count = 0])**

➤ **count**

가

```
>>> re.sub(r'[.,:;]', '', 'a:b;c, d.') #
```

```
'abc d'
```

```
>>> re.sub(r'\W', '', 'a:b;c, d.') # , ,
```

```
#
```

```
'abcd'
```

➤ `{\n}, {g\<n>}` :      `n`  
                                ▪

```
>>> p = re.compile('section{ ( [^}]* ) }', re.VERBOSE)
# VERBOSE
.

>>> p.sub(r'subsection{\1}', 'section{First}
section{second}')
'subsection{First} subsection{second}'
```

➤ `{\g<name>}` :      `n`

▪

```
>>> p = re.compile('section{ (?P<name> [^}]* ) }', re.VERBOSE)
>>> p.sub(r'subsection{\g<name>}', 'section{First}')
'subsection{First}'
```

## ➤ sub

```
>>> def hexrepl( match ):
    value = int( match.group() )
    return hex(value)

>>> p = re.compile(r'\d+')
>>> p.sub(hexrepl, 'Call 65490 for printing, 49152 for user code.')
'Call 0xffd2 for printing, 0xc000 for user code.'
```





➤ 6 - 7

➤ 6 7

```
>>> p = re.compile('(\d{6})-?(\d{7})')
>>> p.search('123456-1234567').groups()
('123456', '1234567')
>>> p.search('1234561234567').groups()
('123456', '1234567')
```

# :HTML

# URL

- HTML

<a href=..>

```
p = re.compile('''  
href=([^\"]\S+?)[\s>] |      # href=http://www.python.or.kr  
href="([^\"]*?)" |          # href="http://www.python.or.kr"  
href='([^\']*?)'            # href='http://www.python.or.kr'  
''', re.I | re.X)
```

```
import re

s = '''<a href="http://www.daum.net">daum</a>
<a href='http://www.naver.com'>naver</a>
<a href=http://www.chosun.com>chosun</a>
<a href=http://www.chosun.com class>chosun</a>
<a href="http://job.daum.net/ " target="new">
<a href="http://go.daum.net/bin/go.cgi?relative=1&url=/Mail-
bin/login_f.cgi%3Ferror%3Dlogin" class="tls0">
'''

p =
re.compile(''href=([^\"]\S+?)[\s>]|href="([^\"]*?)"|href='([^\']*?)\''',
re.I)

pos = 0
while 1:
    match = p.search(s, pos)
    if match:
        url = match.groups()
        pos = match.end()
        print filter(None, url)[0]
    else:
        break
```

●  
  
`http://www.daum.net`

`http://www.naver.com`

`http://www.chosun.com`

`http://www.chosun.com`

`http://job.daum.net/`

`http://go.daum.net/bin/go.cgi?relative=1&url=/Mail-bin/login_f.cgi%3Ferror%3Dlogin`