

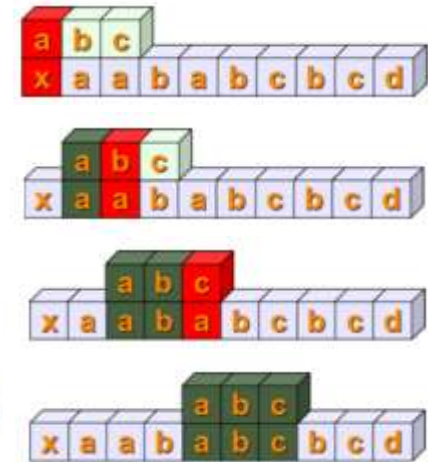
Regular Expressions

Match vs. Search

'abc' in 'xaababcbcd'

`search()` scans through the string to see if any substring matches

```
>>> import re
>>> s = re.search('cat', "A cat and a rat can't be friends.")
>>> print s
<_sre.SRE_Match object at 0x000000000306AE68>
>>> s1 = re.search('cow', "A cat and a rat can't be friends.")
>>> print s1
None
```



`match()` determines if pattern matches at the beginning of a string

```
>>> m = re.match('.*cat.*', "A cat and a rat can't be friends.")
>>> print m
<_sre.SRE_Match object at 0x000000000306AE00>
```

`match()/search()` returns
matched object if matched
None, otherwise

String Pattern Matching

- | means “or”: (aa|bb) matches aa or bb

- [abc] matches a, b, or c

- or simply [a-c]

- equivalent to (a|b|c)

- [abc]+ matches a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc, ...

- [^5] matches any char except 5

- [^0-9] matches any char except a digit

- a[bcd]* matches many more than a(bcd)*

Regex Pattern	Strings Matched
at home	at, home
r2d2 c3po	r2d2, c3po
bat bet bit	bat, bet, bit

Regex Pattern

z.[0-9]

[r-u][env-y]
[us]

[^aeiou]

[^\t\n]

["-a]

String Pattern Matching

- ▶ *****: repeating 0 or more times
 - **ab*d** matches **ad**, **abd**, **abbd**, ...
 - **a(bcd)*d** matches **ad**, **abcdd**, **abcdbcdd**, ...
- ▶ **+**: repeating 1 or more times
 - **ab+d** matches **abd**, **abbd**, ...
 - **a(bcd)+d** matches **abcdd**, **abcdbcdd**, ...
- ▶ **?**: 0 or 1 times
- ▶ **{n}**: n times
- ▶ **{m, n}**: from m to n times

Regex Pattern

`[dn]ot?`

`0?[1-9]`

`[0-9]{15,16}`

`</?[^>]+>`

`[KQRBNP][a-h][1-8]-
[a-h][1-8]`

Special Characters

- characters with special meanings
 - `.` `^` `$` `*` `+` `?` `{` `}` `[` `]` `\` `|` `(` `)`
- `\(ab\)` matches `(ab)`
- `.` matches any single character
 - `.*` matches any string
- `\\` matches `\`
- `^` matches the beginning of a line or string
 - not the `^` inside char-classes `[^...]`
- `$` matches the end of a line or string
 - `a[bcd]*b$` does not match string `abcbdb`

Q: What character class is described by

`[a-zA-Z_][a-zA-Z_0-9]*` ?

Q: regexpr to describe Korean post code?

Regex Pattern

`^From`

`/bin/tcsh$`

`^Subject: hi$`

Predefined Character Classes

- `\d` Matches any decimal digit; equivalent to the set `[0-9]`.
- `\D` The complement of `\d`. It matches any non-digit character; equivalent to the set `[^0-9]`.
- `\s` Matches any whitespace character; equivalent to `[\t\n\r\f\v]`.
- `\S` The complement of `\s`. It matches any non-whitespace character; equiv. to `[^ \t\n\r\f\v]`.
- `\w` Matches any alphanumeric character; equivalent to `[a-zA-Z0-9_]`.
With LOCALE, it will match the set `[a-zA-Z0-9_]` plus characters defined as letters for the current locale.
- `\W` Matches the complement of `\w`.
- `\b` Matches the empty string, but only at the start or end of a word.
- `\B` Matches the empty string, but not at the start or end of a word.

`[\s,.]` matches any white spaces, `,`, `"`, or `"`.

`\b` means word-boundary (zero-length):

`\b\w+\b` matches a single word (actually `\b\w+` is enough)

Regex Pattern	Strings Matched
<code>\w+-\d+</code>	Alphanumeric string and number separated by a hyphen
<code>[A-Za-z]\w*</code>	Alphabetic first character; additional characters (if present) can be alphanumeric (almost equivalent to the set of valid Python identifiers [see exercises])
<code>\d{3}-\d{3}-\d{4}</code>	American-format telephone numbers with an area code prefix, as in 800-555-1212
<code>\w+@\w+\.</code> <i>com</i>	Simple e-mail addresses of the form <i>XXX@YYY.com</i>

Regex Pattern	Strings Matched
<code>\d+(\.\d*)?</code>	Strings representing simple floating-point numbers; that is, any number of digits followed optionally by a single decimal point and zero or more numeric digits, as in "0.004," "2," "75.," etc.
<code>(Mr?s?\.\s)?[A-Z][a-z]* [A-Za-z-]+</code>	First name and last name, with a restricted first name (must start with uppercase; lowercase only for remaining letters, if any), the full name, prepended by an optional title of "Mr.," "Mrs.," "Ms.," or "M.," and a flexible last name, allowing for multiple words, dashes, and uppercase letters

Matching is Greedy

- by “matching” we mean matching the beginning portion of a string
 - $a(bc)^+$ matches the underlined part in abcbcd
- greedy search with backtracking
- $a(bcd)^*b$ matches abcbd, abcbcd, abcd
- try match the pattern $a[bcd]^*b$ with string abcbd

1	a	The a in the RE matches.
2	abcbcd	The engine matches $[bcd]^*$, going as far as it can, which is to the end of the string.
3	<i>Failure</i>	The engine tries to match b, but the current position is at the end of the string, so it fails.
4	abcb	Back up, so that $[bcd]^*$ matches one less character.
5	<i>Failure</i>	Try b again, but the current position is at the last character, which is a "d".
6	abc	Back up again, so that $[bcd]^*$ is only matching "bc".
6	abcb	Try b again. This time but the character at the current position is "b", so it succeeds.

Performing Match and Search

```
>>> import re

>>> re.match('[a-z]+', '')
None

>>> p = re.compile('[a-z]+')
>>> p
<_sre.SRE_Pattern object at 80c3c28>

>>> p.match('')
>>> print p.match('')
None

>>> m = p.match('tempo')
>>> print m
<_sre.SRE_Match object at 80c4f68>
```

```
>>> print p.match('::: message')
None
>>> m = p.search('::: message')
>>> print m
<re.MatchObject instance at 80c9650>

>>> m.group()
'message'
>>> m.span()
(4, 11)
```

compiled version is faster
for repeated use

Grouping

```
>>> import re
>>> patt = '\w+@(\w+\.)?\w+\.com'
>>> re.match(patt, 'nobody@www.xxx.com').group()
'nobody@www.xxx.com'
>>>
>>> group_patt = '(\w+)@(\w+\.)?\w+\.com'
>>> p = re.compile(group_patt)
>>> m = p.match('nobody@www.xxx.com')
>>> m.group()    # entire match
'nobody@www.xxx.com'
>>> m.group(0)   # same as group()
'nobody@www.xxx.com'
>>> m.group(1)   # subgroup 1
'nobody'
>>> m.group(2)   # subgroup 2
'www.'
>>> m.groups()   # all subgroups
('nobody', 'www.')
```

Common Match Object Methods (see documentation for others)

<code>group(num=0)</code>	Return entire match (or specific subgroup <i>num</i>)
<code>groups</code> (<i>default=None</i>)	Return all matching subgroups in a tuple (empty if there aren't any)
<code>groupdict</code> (<i>default=None</i>)	Return dict containing all matching named subgroups with the names as the keys (empty if there weren't any)

findall() vs. finditer()

- `findall()` returns a **list** of all substrings that matches
- `finditer()` returns an **iterator** of **matched objects**

```
>>> p = re.compile('\d+')
>>> s = '12 drummers, 11 pipers, 10 lords'
>>> p.findall()
['12', '11', '10']

>>> iterator = p.finditer(s)
>>> iterator
<callable-iterator object at 0x4018...>
>>> for match in iterator:
...     print match.span()
...
(0, 2)
(13, 15)
(24, 26)
```

```
>>> import re
>>> patt = re.compile(r'(\w+)\s+(\d+)')
>>> s = " I teach cs 399 and cis 500."
>>>
>>> m = patt.search(s)
>>> m.group()
'cs 399'
>>> m.groups()
('cs', '399')
>>>
>>> patt.findall(s)
[('cs', '399'), ('cis', '500')]
>>>
>>> for m in p.finditer(s):
...     print m.group(), m.groups()
...
cs 399 ('cs', '399')
cis 500 ('cis', '500')
```

Back-referencing Groups

- referring to previous groups by `\1`, `\2`, ...

```
>>> q = re.compile(r'(\b\w+)\s+\1')  
>>> s = "this is the the course"
```

r means "raw"

```
>>> q.findall(s)  
['the']
```

typo detector :)

```
>>> q.search(s).group()  
'the the'
```

```
>>> p = re.compile('(a(b)c)d')  
>>> m = p.match('abcd')  
>>> m.group(0)  
'abcd'  
>>> m.group(1)  
'abc'  
>>> m.group(2)  
'b'
```


Non-greedy Quantifier

- default matching is greedy
- use non-greedy quantifiers ?

```
>>> s = '<html><head><title>Title</title>'
```

```
>>> print re.match('<.*>', s).group()  
<html><head><title>Title</title>
```

```
>>> print re.match('<.*?>', s).group()  
<html>
```

```
>>> p = re.compile('<a href=(.*)>')  
>>> p.match("<a href=\"index.html\">back</a>").group(1)  
"index.html">back</a>
```

```
>>> p = re.compile('<a href=(.*?)>')  
>>> p.match("<a href=\"index.html\">back</a>").group(1)  
"index.html"
```

Common Regular Expression Attribute

re Module Function Only

<code>compile(pattern, flags=0)</code>	Compile regex <i>pattern</i> with any optional <i>flags</i> and return a regex object
--	---

re Module Functions and Regex Object Methods

<code>match(pattern, string, flags=0)</code>	Attempt to match <i>pattern</i> to <i>string</i> with optional <i>flags</i> ; return match object on success, None on failure
--	---

<code>search(pattern, string, flags=0)</code>	Search for first occurrence of <i>pattern</i> within <i>string</i> with optional <i>flags</i> ; return match object on success, None on failure
---	---

<code>findall(pattern, string[, flags])^a</code>	Look for all (non-overlapping) occurrences of <i>pattern</i> in <i>string</i> ; return a list of matches
--	--

<code>finditer(pattern, string[, flags])^b</code>	Same as <code>findall()</code> , except returns an iterator instead of a list; for each match, the iterator returns a match object
---	--

<code>split(pattern, string, max=0)^c</code>	Split <i>string</i> into a list according to regex <i>pattern</i> delimiter and return list of successful matches, splitting at most <i>max</i> times (split all occurrences is the default)
--	--

<code>sub(pattern, repl, string, count=0)^c</code>	Replace all occurrences of the regex <i>pattern</i> in <i>string</i> with <i>repl</i> , substituting all occurrences unless <i>count</i> provided (see also <code>subn()</code> , which, in addition, returns the number of substitutions made)
--	---

<code>purge()</code>	Purge cache of implicitly compiled regex patterns
----------------------	---

Common Module Attributes (flags for most regex functions)

<code>re.I, re.IGNORECASE</code>	Case-insensitive matching
----------------------------------	---------------------------

<code>re.L, re.LOCALE</code>	Matches via <code>\w, \W, \b, \B, \s, \S</code> depends on locale
------------------------------	---

<code>re.M, re.MULTILINE</code>	Respectively causes <code>^</code> and <code>\$</code> to match the beginning and end of each line in target string rather than strictly the beginning and end of the entire string itself
---------------------------------	--

<code>re.S, re.DOTALL</code>	The <code>.</code> normally matches any single character except <code>\n</code> ; this flag says <code>.</code> should match them, too
------------------------------	--

<code>re.X, re.VERBOSE</code>	All whitespace plus <code>#</code> (and all text after it on a single line) are ignored unless in a character class or backslash-escaped, allowing comments and improving readability
-------------------------------	---