

Regular Expressions in programming

CSE 307 – Principles of Programming Languages

Stony Brook University

<http://www.cs.stonybrook.edu/~cse307>

What are *Regular Expressions*?

- Formal language representing a text pattern interpreted by a regular expression processor
 - Used for **matching**, **searching** and **replacing** text
 - There are no variables and you cannot do mathematical operations (such as: you cannot add **1+1**) – it is not a programming language
 - Frequently you will hear them called *regex* or *RE* for short (or pluralized "*regexes*")

What are Regular Expressions?

- Usage examples:
 - Test if a phone number has the correct number of digits
 - Test if an email address has the correct format
 - Test if a Social Security Number is in the correct format
 - Search a text for words that contain digits
 - Find duplicate words in a text
 - Replace all occurrences of "*Bob*" and "*Bobby*" with "*Robert*"
 - Count the number of times "*science*" is preceded by "*computer*" or "*information*"
 - Convert a tab indentations file with spaces indentations

What are Regular Expressions?

- But what is "*Matches*"?
 - a text *matches* a regular expression if it is correctly described by the regex

```
>>> m = re.match(r"(\w+) (\w+)", "Isaac Newton, physicist")
>>> m
<re.Match object; span=(0, 12), match='Isaac Newton'>
```

```
>>> m.group(0) # The entire match
'Isaac Newton'
```

```
>>> m.group(1) # The first parenthesized subgroup.
'Isaac'
```

```
>>> m.group(2) # The second parenthesized subgroup.
'Newton'
```

History of Regular Expressions

- 1943: Warren McCulloch and Walter Pitts developed models of how the nervous system works
- 1956: Steven Kleene described these models with an algebra called "*regular sets*" and created a notation to express them called "*regular expressions*"
- 1968: Ken Thompson implements regular expressions in **ed**, a Unix text editor
 - Example: **g/Regular Expression/p**
 - meaning Global Regular Expression Print (**grep**)
 - g = global / whole file; p = print

History of Regular Expressions

- grep evolved into egrep
 - but broke backward compatibility
- Therefore, in 1986, everyone came together and defined POSIX (Portable Operating Systems Interface)
 - Basic Regular Expressions (BREs)
 - Extended Regular Expressions (EREs)
- 1986: Henry Spencer releases the **regex** library in C
 - Many incorporated it in other languages and tools
- 1987: Larry Wall released Perl
 - Used Spencer's regex library
 - Added powerful features
 - Everybody wanted to have it in their languages: *Perl Compatible RE (PCRE)* library, Java, Javascript, C#/VB/.NET, MySQL, PHP, Python, Ruby

Regular Expressions Engines

- Main versions / standards:
 - PCRE
 - POSIX BRE
 - POSIX ERE
- Very subtle differences
 - Mainly older UNIX tools that use POSIX BRE for compatibility reasons
- In use:
 - Unix (POSIX BRE, POSIX ERE)
 - PHP (PCRE)
 - Apache (v1: POSIX ERE, v2: PCRE)
 - MySQL (POSIX ERE)
- Each of these languages is improving, so check their manuals

Python Regular Expressions

- <https://docs.python.org/3/library/re.html>

- It is more powerful than String splits:

```
>>> "ab bc cd".split()
['ab', 'bc', 'cd']
```

- Import the re module:

```
import re
```

```
>>> re.split(" ", "ab bc cd")
['ab', 'bc', 'cd']
```

```
>>> re.split("\d", "ab1bc4cd")
['ab', 'bc', 'cd']
```

```
>>> re.split("\d*", "ab13bc44cd443gg")
['', 'a', 'b', '', 'b', 'c', '', 'c', 'd',
'', 'g', 'g', '']
```


Python Regular Expressions

```
>>> re.split("\d+", "ab13bc44cd443gg")  
['ab', 'bc', 'cd', 'gg']
```

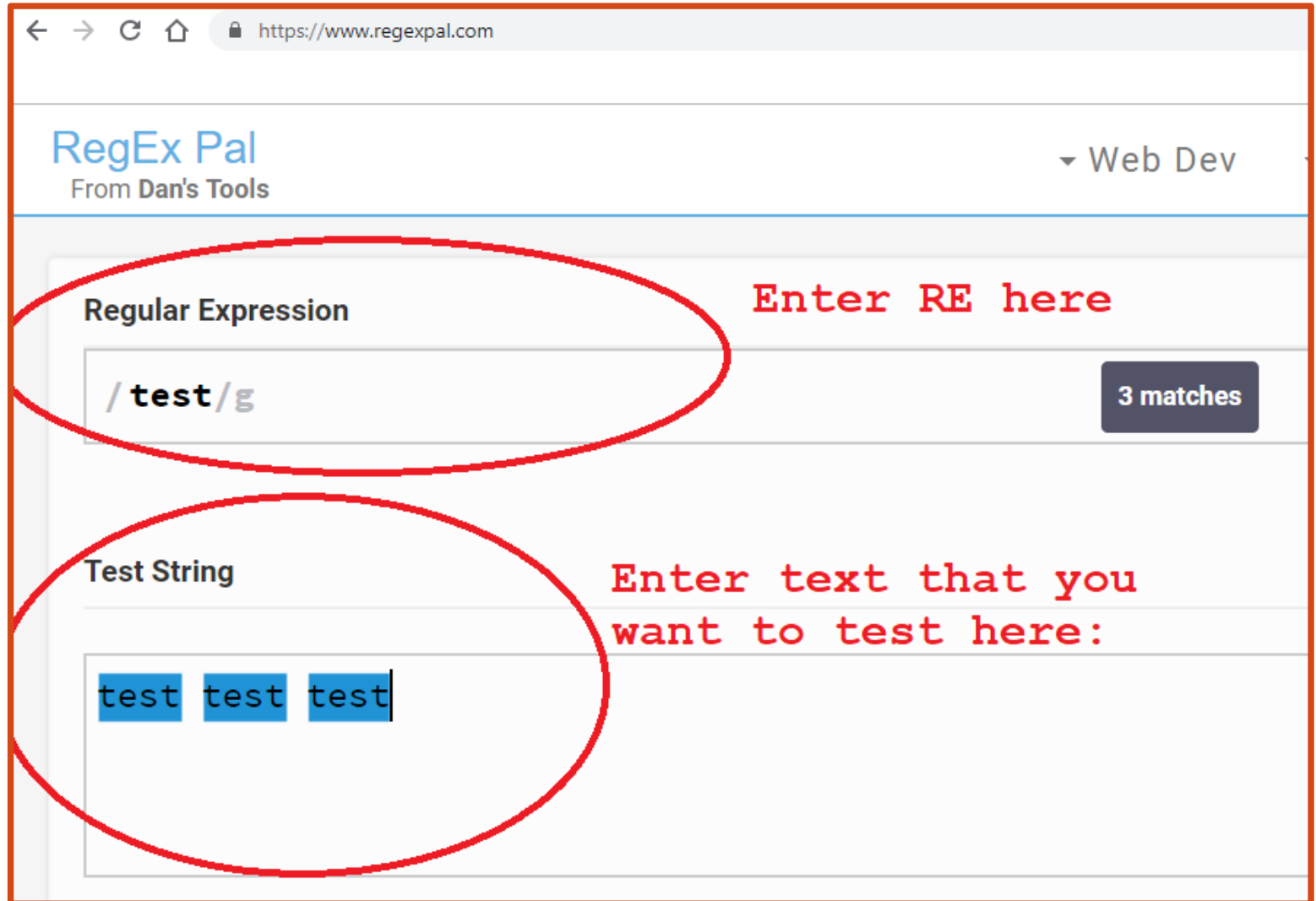
```
>>> m = re.search('(?!<=abc)def', 'abcdef')
```

```
>>> m
```

```
<re.Match object; span=(3, 6), match='def'>
```

Online Regular Expressions

- <https://regexpal.com>



The screenshot shows the RegEx Pal website interface. The browser address bar displays `https://www.regexpal.com`. The page title is "RegEx Pal" with the subtitle "From Dan's Tools" and a "Web Dev" dropdown menu. The main content area is divided into two sections: "Regular Expression" and "Test String".

In the "Regular Expression" section, the input field contains `/test/g`, which is circled in red. To the right of this field, the text "Enter RE here" is displayed in red. A dark blue button on the right indicates "3 matches".

In the "Test String" section, the input field contains the text "test test test", which is also circled in red. To the right of this field, the text "Enter text that you want to test here:" is displayed in red.

Regular Expressions

- Strings:
 - "car" matches "car"
 - "car" also matches the first three letters in "cartoon"
 - "car" does not match "c_a_r"
- Similar to search in a word processor
- Case-sensitive (by default): "car" does not match "Car"
- Metacharacters:
 - Have a special meaning
 - Like mathematical operators
 - Transform char sequences into powerful patterns
 - Only a few characters to learn: \ . * + - { } [] () ^ \$ | ? : ! =
 - May have multiple meanings
 - Depend on the context in which they are used
 - Variation between regex engines

The *wildcard character*

- Like in card games: one card can replace any other card on the pattern

Metacharacter	Meaning
.	Any character except newline

- Examples:
 - "h.t" matches "hat", "hot", "heat"
 - ".a.a.a" matches "banana", "papaya"
 - "h.t" does not match "'heat" or "Hot"
- Common mistake:
 - "9.00" matches "9.00", but it also match "9500", "9-00"
- **We should write regular expressions to match what we want and ONLY what we want (We don't want to be overly permissive, we don't want false positives, we want the regular expression to match what we are not looking for)**

Escaping Metacharacter

- Allow the use of metacharacters **as characters**:

- "\." matches "."

Metacharacter	Meaning
\	Escape the next character

- "9\.00" matches only "9.00", but not "9500" or "9-00"
- Match a backslash by escaping it with a backslash:
 - "\\\" matches only "\"
 - "C:\\Users\\Paul" matches "C:\Users\Paul"
- Only for metacharacters
 - literal characters should never be escaped because it gives them meaning, e.g., r"\n"
 - Sometimes we want both meanings
 - Example: we want to match files of the name: "1_report.txt", "2_report.txt", ...
 - "._report\\.txt" uses the first . as a wildcard and the second \. as the period itself

Other special characters

- Tabs: `\t`
- Line returns: `\r` (*line return*), `\n` (*newline*), `\r\n`
- Unicode codes: `\u00A9`
- ASCII codes: `\x00A9`

Character sets

Metacharacter	Meaning
[Begin character set
]	End character set

- Matches any of the characters inside the set
 - But only one character
 - Order of characters does not matter
 - Examples:
 - "[aeiou]" matches a single vowel, such as: "a" or "e"
 - "gr[ae]y" matches "gray" or "grey"
 - "gr[ae]t" does not match "great"

Character ranges

- `[a-z] = [abcdefghijklmnopqrstuvwxyz]`
 - Range metacharacter - is only a character range when it is inside a character set, a dash line otherwise
 - represent all characters between two characters
 - Examples:
 - `[0-9]`
 - `[A-Za-z]`
 - `[0-9A-Za-z]`
 - `[0-9][0-9][0-9]-[0-9][0-9][0-9]-[0-9][0-9][0-9][0-9]` matches phone "631-632-9820"
 - `[0-9][0-9][0-9][0-9][0-9]` matches zip code "90210"
 - `[A-Z0-9][A-Z0-9][A-Z0-9] [A-Z0-9][A-Z0-9][A-Z0-9]` matches Canadian zip codes, such as, "VC5 B6T"
 - Caution:
 - What is `[50-99]`?
 - It is not $\{50, 51, \dots, 99\}$
 - It is same with `[0-9]`: the set contains already 5 and 9

Negative Character sets

Metacharacter	Meaning
<code>^</code>	Negate a character set

- Caret (`^`) = not one of several characters
 - Add `^` as the first character inside a character set
 - Still represents one character
 - Examples:
 - `[^aeiou]` matches any one character that is not a lower case vowel
 - `[^aeiouAEIOU]` matches any one character that is not a vowel (non-vowel)
 - `[^a-zA-Z]` matches any one character that is not a letter
 - `see[^mn]` matches "seek" and "sees", but not "seem" or "seen"
 - `see[^mn]` matches "see " because space matches `[^mn]`
 - `see[^mn]` does not match "see" because there is no more character after see

Metacharacters

- Metacharacters inside Character sets are already escaped:
 - Do not need to escape them again
 - Examples:
 - `h[o.]t` matches "hot" and "h.t"
 - Exceptions: metacharacters that have to do with character sets: `]-^\`
 - Examples:
 - `[[\]]` matches "[" or "]"
 - `var[[()[]0-9]()\]]` matches "var()" or "var[]"
- Exception to exception: `"10[-/]10"` matches "10-10" or "10/10"
 - - does not need to be escaped because it is not a range

Shorthand character sets

Shorthand	Meaning	Equivalent
<code>\d</code>	Digit	<code>[0-9]</code>
<code>\w</code>	Word character	<code>[a-zA-z0-9_]</code>
<code>\s</code>	Whitespace	<code>[\t\n\r]</code>
<code>\D</code>	Not digit	<code>[^0-9]</code>
<code>\W</code>	Not word character	<code>[^a-zA-z0-9_]</code>
<code>\S</code>	Not white space	<code>[^ \t\n\r]</code>

- Underscore (`_`) is a word character
- Hyphen (`-`) is not a word character
 - `"\d\d\d"` matches `"123"`
 - `"\w\w\w"` matches `"123"` and `"ABC"` and `"1_A"`
 - `"\w\s\w\w"` matches `"I am"`, but not `"Am I"`
 - `"[^d]"` matches `"a"`
 - `"[^d\w]"` is not the same with `"[\D\W]"` (accepts `"a"`)

Introduced in Perl
Not in many Unix tools

POSIX Bracket Expressions

POSIX	Description	ASCII	Unicode	Shorthand	Java
<code>[a-zA-Z0-9]</code>	Alphanumeric characters	<code>[a-zA-Z0-9]</code>	<code>[\p{L}\p{N}]</code> <code>[\p{Nd}]</code>		<code>\p{Alnum}</code>
<code>[a-zA-Z]</code>	Alphabetic characters	<code>[a-zA-Z]</code>	<code>\p{L}\p{N}</code>		<code>\p{Alpha}</code>
<code>[\x00-\x7F]</code>	ASCII characters	<code>[\x00-\x7F]</code>	<code>\p{InBasicLatin}</code>		<code>\p{ASCII}</code>
<code>[\t]</code>	Space and tab	<code>[\t]</code>	<code>[\p{Zs}\t]</code>	<code>\h</code>	<code>\p{Blank}</code>
<code>[\x00-\x1F\x7F]</code>	Control characters	<code>[\x00-\x1F\x7F]</code>	<code>\p{Cc}</code>		<code>\p{Cntrl}</code>
<code>[0-9]</code>	Digits	<code>[0-9]</code>	<code>\p{Nd}</code>	<code>\d</code>	<code>\p{Digit}</code>
<code>[\x21-\x7E]</code>	Visible characters (anything except spaces and control characters)	<code>[\x21-\x7E]</code>	<code>[\p{Z}\p{C}]</code>		<code>\p{Graph}</code>
<code>[a-z]</code>	Lowercase letters	<code>[a-z]</code>	<code>\p{L}</code>	<code>\l</code>	<code>\p{Lower}</code>
<code>[\x20-\x7E]</code>	Visible characters and spaces (anything except control characters)	<code>[\x20-\x7E]</code>	<code>\P{C}</code>		<code>\p{Print}</code>
<code>[!"#\$%&'()*+,-./:;<=>?@\[\]^_`{ }~]</code>	Punctuation (and symbols).	<code>[!"#\$%&'()*+,-./:;<=>?@\[\]^_`{ }~]</code>	<code>\p{P}</code>		<code>\p{Punct}</code>
<code>[\t\r\n\v\f]</code>	All whitespace characters, including line breaks	<code>[\t\r\n\v\f]</code>	<code>[\p{Z}\t\r\n\v\f]</code>	<code>\s</code>	<code>\p{Space}</code>
<code>[A-Z]</code>	Uppercase letters	<code>[A-Z]</code>	<code>\p{Lu}</code>	<code>\u</code>	<code>\p{Upper}</code>
<code>[A-Za-z0-9_]</code>	Word characters (letters, numbers and underscores)	<code>[A-Za-z0-9_]</code>	<code>[\p{L}\p{N}]</code> <code>[\p{Nd}\p{Pc}]</code>	<code>\w</code>	<code>\p{Isword}</code>
<code>[A-Fa-f0-9]</code>	Hexadecimal digits	<code>[A-Fa-f0-9]</code>	<code>[A-Fa-f0-9]</code>		<code>\p{XDigit}</code>

Repetition

Metacharacter	Meaning
*	Preceding item zero or more times
+	Preceding item one or more times
?	Preceding item zero or one time

- Examples:
 - `apples*` matches "apple" and "apples" and "applessssssss"
 - `apples+` matches "apples" and "applessssssss"
 - `apples?` matches "apple" and "apples"
 - `\d*` matches "123"
 - `colou?r` matches "color" and "colour"

Quantified Repetition

Metacharacter	Meaning
{	Start quantified repetition of preceding item
}	End quantified repetition of preceding item

- {min, max}
 - min and max must be positive numbers
 - min must always be included
 - min can be 0
 - max is optional
- Syntaxes:
 - `\d{4,8}` matches numbers with 4 to 8 digits
 - `\d{4}` matches numbers with exactly 4 digits
 - `\d{4,}` matches numbers with minimum 4 digits
 - `\d{0,}` is equivalent to `\d*`
 - `\d{1,}` is equivalent to `\d+`

Greedy Expressions

- Standard repetition quantifiers are **greedy**:
 - expressions try to match the longest possible string
 - `\d*` matches the entire string "1234" and not just "123", "1", or "23"
- Lazy expressions:
 - matches as little as possible before giving control to the next expression part
 - `?` makes the preceding quantifier into a lazy quantifier
 - `*?`
 - `+?`
 - `{min,max}?`
 - `??`
 - Example:
 - "apples??" matches "apple" in "apples"

Grouping metacharacters

Metacharacter	Meaning
(Start grouped expression
)	End grouped expression

- Group a large part to apply repetition to it
 - "(abc)*" matches "abc" and "abcabcabc"
 - "(in)?dependent" matches "dependent" and "independent"
- Makes expressions easier to read
- Cannot be used inside character sets

Metacharacters

- `$` Matches the ending position of the string or the position just before a string-ending newline.
 - In line-based tools, it matches the ending position of any line.
 - `[hc]at$` matches "hat" and "cat", but only at the end of the string or line.
- `^` Matches the beginning of a line or string.
- `|` The choice (also known as alternation or set union) operator matches either the expression before or the expression after the operator.
 - For example, `abc|def` matches "abc" or "def".
- `\A` Matches the beginning of a string (but not an internal line).
- `\z` Matches the end of a string (but not an internal line).

Summary: Frequently Used Regular Expressions

Regular Expression	Meaning	Example
x	A character literal	"good" matches "good"
.	Any single character	"good" matches "goo."
(ab cd)	ab or cd	"good" matches "a g"
[abc]	a, b, or c	"good" matches "[ag]"
[^abc]	any character except a, b, or c	"good" matches "[^ac]"
[a-z]	a through z	"good" matches <u>[a-i]oo[a-d]</u>
[^a-z]	any character except a through z	"good" matches <u>goo[^i-x]</u>
\d	a digit, same as [0-9]	"good3" matches <u>"good\d"</u>
\D	a non-digit	"good" matches <u>"\D\Dod"</u>
\w	a word character	"good3" matches <u>"goo\w\w"</u>
\W	a non-word character	<u>\$good</u> matches <u>"\Wgood"</u>
\s	a whitespace character	<u>"good 2"</u> matches <u>"good\s2"</u>
\S	a non-whitespace char	"good" matches <u>"\Sood"</u>
p*	zero or more occurrences of pattern p	"good" matches <u>"a*"</u> <u>bbb</u> matches <u>"a*"</u>
p+	one or more occurrences of pattern p	"good" matches <u>"o+"</u> <u>bbb</u> matches <u>"b+"</u>
p?	zero or one occurrence of pattern p	"good" matches <u>"good?"</u> <u>bbb</u> matches <u>"b?"</u>
p{n}	exactly n occurrences of pattern p	<u>aaa</u> matches <u>"a{3}"</u>
p{n,}	at least n occurrences of pattern p	<u>good</u> does not match <u>"go{2}d"</u> <u>good</u> matches <u>"go{2,}d"</u>
p{n,m}	between n and m occurrences (inclusive)	<u>good</u> does not match <u>"g{1,}"</u> <u>aa</u> matches <u>"a{1,9}"</u> <u>bb</u> does not match <u>"b{2,9}"</u>

Python `match` and `search` Functions

- `re.match(r, s)` returns a match object if the regex `r` matches at the start of string `s`

```
import re
regex = "\d{3}-\d{2}-\d{4}"
ssn = input("Enter SSN: ")
match1 = re.match(regex, ssn)
if match1 != None:
    print(ssn, " is a valid SSN")
    print("start position of the matched text is "
          + str(match1.start()))
    print("start and end position of the matched text is "
          + str(match1.span()))
else:
    print(ssn, " is not a valid SSN")
```

```
Enter SSN: 123-12-1234 more text
123-12-1234 more text is a valid SSN
start position of the matched text is 0
start and end position of the matched text is (0, 11)
```

Python `match` and `search` Functions

- Invoking `re.match` returns a match object if the string matches the regex pattern at the start of the string.
 - Otherwise, it returns **None**.
- The program checks whether if there is a match.
 - If so, it invokes the match object's `start()` method to return the start position of the matched text in the string (line 10) and the `span()` method to return the start and end position of the matched text in a tuple (line 11).

Python `match` and `search` Functions

- `re.search(r, s)` returns a match object if the regex `r` matches anywhere in string `s`

```
import re
regex = "\d{3}-\d{2}-\d{4}"
text = input("Enter a text: ")
match1 = re.search(regex, text)
if match1 != None:
    print(text, " contains a valid SSN")
    print("start position of the matched text is "
          + str(match1.start()))
    print("start and end position of the matched text is "
          + str(match1.span()))
else:
    print(ssn, " does not contain a valid SSN")
```

```
Enter a text: The ssn for Smith is 343-34-3490
The ssn for Smith is 343-34-3490 contains a SSN
start position of the matched text is 21
start and end position of the matched text is (21, 32)
```

Flags

- For the functions in the **re** module, an optional flag parameter can be used to specify additional constraints
- For example, in the following statement

```
re.search("a{3}", "AaaBe", re.IGNORECASE)
```

The string "AaaBe" matches the pattern `a{3}` case-insensitive

Findall

- **findall(pattern, string [, flags])** return a list of strings giving all nonoverlapping matches of pattern in string. If there are any groups in patterns, returns a list of groups, and a list of tuples if the pattern has more than one group

```
>>> re.findall('<(.*?)>', '<spam> /<ham><eggs>')  
['spam', 'ham', 'eggs']
```

```
>>> re.findall('<(.*?)>/?<(.*?)>',  
              '<spam>/<ham> ... <eggs><cheese>')  
[('spam', 'ham'), ('eggs', 'cheese')]
```

Findall

- **sub(pattern, repl, string [, count, flags])** returns the string obtained by replacing the (first count) leftmost nonoverlapping occurrences of pattern (a string or a pattern object) in string by **repl** (which may be a string with backslash escapes that may back-reference a matched group, or a function that is passed a single match object and returns the replacement string).
- **compile(pattern [, flags])** compiles a regular expression pattern string into a regular expression pattern object, for later matching.

Groups

- Groups: extract substrings matched by REs in '()' parts
 - (R) Matches any regular expression inside (), and delimits a group (retains matched substring)
 - \N Matches the contents of the group of the same number N: '(.) \1' matches "42 42"

```
import re
patt = re.compile("A(.)B(.)C(.)") # saves 3 substrings
mobj = patt.match("A0B1C2") # each '()' is a group, 1..n
print(mobj.group(1), mobj.group(2), mobj.group(3))
patt = re.compile("A(*)B(*)C(*)") # saves 3 substrings
mobj = patt.match("A000B111C222") # groups() gives all groups
print(mobj.groups())
print(re.search("(A|X)(B|Y)(C|Z)D", "..AYCD..").groups())
print(re.search("(?P<a>A|X)(?P<b>B|Y)(?P<c>C|Z)D",
    "..AYCD..").groupdict())
patt = re.compile(r"[\t ]*#\s*define\s*([a-z0-9_]*)\s*(.*)")
mobj = patt.search("# define spam 1 + 2 + 3") # parts of C #define
print(mobj.groups()) # \s is whitespace
```

Groups

```
python re-groups.py
```

```
0 1 2
```

```
('000', '111', '222')
```

```
('A', 'Y', 'C')
```

```
{'a': 'A', 'c': 'C', 'b': 'Y'}
```

```
('spam', '1 + 2 + 3')
```

Groups

- When a match or search function or method is successful, you get back a match object
 - **group(g)** **group(g1, g2, ...)** Return the substring that matched a parenthesized group (or groups) in the pattern. Accept group numbers or names. Group numbers start at 1; group 0 is the entire string matched by the pattern. Returns a tuple when passed multiple group numbers, and group number defaults to 0 if omitted
 - **groups()** Returns a tuple of all groups' substrings of the match (for group numbers 1 and higher).
 - **start([group])** **end([group])** Indices of the start and end of the substring matched by group (or the entire matched string, if no group is passed).
 - **span([group])** Returns the two-item tuple: **(start(group), end(group))**