

Discussion #4

Name:

Regular Expressions

Here's a complete list of metacharacters:

. ^ \$ * + ? { } [] \ | ()

Some reminders on what each can do (this is not exhaustive):

- | | |
|---|---|
| "^" matches the position at the beginning of string (unless used for negation "[^"]) | "["]" match any one of the characters inside, accepts a range, e.g., "[a-c]" |
| "\$" matches the position at the end of string character. | "(")" used to create a sub-expression |
| "?" match preceding literal or sub-expression 0 or 1 times. When following "+" or "*" results in non-greedy matching. | "\d" match any <i>digit</i> character. "\D" is the complement. |
| "+" match preceding literal or sub-expression <i>one</i> or more times. | "\w" match any <i>word</i> character (letters, digits, underscore). "\W" is the complement. |
| "*" match preceding literal or sub-expression <i>zero</i> or more times | "\s" match any <i>whitespace</i> character including tabs and newlines. \S is the complement. |
| "." match any character except new line. | "\b" match boundary between words |

Some useful `re` package functions:

- | | |
|--|--|
| <code>re.split(pattern, string)</code> split the string at substrings that match the pattern. Returns a list. | <code>re.sub(pattern, replace, string)</code> apply the pattern to string replacing matching substrings with <code>replace</code> . Returns a string. |
|--|--|

Reading Regex

1. Given the text,

```
<record> Joseph Gonzalez <jegonzal@berkeley.edu> Faculty </record>
<record> Jake Soloff <jake_soloff@berkeley.edu> TA </record>
```

Which of the following matches exactly to the email addresses (including angle brackets)?

- (a) `<.*@.*>`
 (b) `<\w*@.*?>`
 (c) `<[^>]*>`
2. Which strings contain a match for the following regular expression, `abc?&`
- (a) Know your abcs
 (b) Did you say abc?
 (c) Hi ab
3. For each pattern specify the starting and ending position of the first match in the string.

	abcdefg	abcs!	ab abc	abc, 123
<code>abc*</code>	1-3			
<code>[^\s]+</code>				
<code>ab.*c</code>				
<code>[a-z1,9]+</code>				

Writing Regex

- 4.
- (a) Write a regular expression that matches a string that contains only lowercase letters and numbers (including empty string).
- (b) Given `text = "123 Fake Street"`, use methods in RE module to abbreviate "`Street`" as "`St.`". The result should look like "`123 Fake St.`".
- (c) Given `text2 = "October 10, November 11, December 12, January 1"`, use methods in RE module to extract all the numbers in the string. The result should look like `["10", "11", "12", "1"]`.