



Utrecht
Bioinformatics
Center

Regular Expression

Day 4b – Introduction to Python



Regular Expression

- Wikipedia:
- “a regular expression (regex) is a sequence of characters that forms a search pattern.”
- What can we do with them?
 - 1. Find matches in target (strings)
 - 2. Substitute patterns

Regular Expression



Anchors		Sample Patterns	
^	Start of line +	([A-Za-z0-9-]+)	Letters, numbers and hyphens
\A	Start of string +	(\d{1,2}\V\d{1,2}\V\d{4})	Date (e.g. 21/3/2006)
\$	End of line +	([^\s]+(?:\.(jpg gif png))\.\2)	jpg, gif or png image
\Z	End of string +	(^[1-9]{1}\$ ^[1-4]{1}[0-9]{1}\$ ^50\$)	Any number from 1 to 50 inclusive
\b	Word boundary +	(#?([A-Fa-f0-9]){3}((([A-Fa-f0-9]){3})?))	Valid hexadecimal colour code
\B	Not word boundary +	((?=\.*\d)(?=\.*[a-z])(?=\.*[A-Z]).{8,15})	8 to 15 character string with at least one upper case letter, one lower case letter, and one digit (useful for passwords).
\<	Start of word	(\w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6})	Email addresses
\>	End of word	(\<(/?[^\>]+)\>)	HTML Tags
Character Classes			
\c	Control character	Note <i>These patterns are intended for reference purposes and have not been extensively tested. Please use with caution and test thoroughly before use.</i>	
\s	White space		
\S	Not white space		
\d	Digit		
\D	Not digit		
\w	Word		
\W	Not word		
\xhh	Hexadecimal character hh		
\Oxxx	Octal character xxx		
POSIX Character Classes			
[[:upper:]]	Upper case letters	Quantifiers	
[[:lower:]]	Lower case letters		
[[:alpha:]]	All letters		
[[:alnum:]]	Digits and letters		
[[:digit:]]	Digits		
[[:xdigit:]]	Hexadecimal digits		
[[:punct:]]	Punctuation		
[[:blank:]]	Space and tab		
[[:space:]]	Blank characters	Ranges	
[[:cntrl:]]	Control characters		
[[:graph:]]	Printed characters		
[[:print:]]	Printed characters and spaces		
[[:word:]]	Digits, letters and underscore		
[[:upper:]]	Upper case letters		
[[:lower:]]	Lower case letters		
[[:alpha:]]	All letters		
[[:alnum:]]	Digits and letters		
[[:digit:]]	Digits		
[[:xdigit:]]	Hexadecimal digits		
[[:punct:]]	Punctuation		
[[:blank:]]	Space and tab		
[[:space:]]	Blank characters		
[[:cntrl:]]	Control characters		
[[:graph:]]	Printed characters		
[[:print:]]	Printed characters and spaces		
[[:word:]]	Digits, letters and underscore		
Assertions		Special Characters	
?=	Lookahead assertion +		
?!	Negative lookahead +		
?<=	Lookbehind assertion +		
?!= or ?<!	Negative lookbehind +		
?>	Once-only Subexpression		
?()	Condition [if then]		
?()	Condition [if then else]		
?#	Comment		
Note		Items marked + should work in most regular expression implementations.	

RegEx - Intro

- Keep it simple: Operators
- Example: Find EcoRI restriction enzyme site in sequence.

```
EcoRI = "GAATTC"  
DNAseq = "TGCATAGCGAATTCGGACGT"  
EcoRI in DNAseq  
True
```

RegEx - Intro

- Find Eco13kl restriction site CCNGG

```
Eco13kl = "CCNGG"  
DNAseq = "CCTGGAGCCCAGGGGACGT"  
Eco13kl in DNAseq  
False
```

- CCNGG = CCAGG or CCTTGG or CCCCGG or CCGGGG

```
Eco13kl=["CCAGG","CCTGG","CCCGG","CCGGG"]  
DNAseq = "CCTGGAGCCCAGGGGACGT"  
Eco13kl[0] in DNAseq      True  
Eco13kl[1] in DNAseq      True  
Eco13kl[2] in DNAseq      False  
Eco13kl[3] in DNAseq      False
```

RegEx - module

- For complicated patterns use Regex module “re”

```
import re
```

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

RegEx – re.findall

- `re.findall(pattern, string, flags=0)`
- Return all non-overlapping matches of pattern in string, as a list of strings.
- Pattern → pattern to search
- String → string to search in
- Flags → optional flags/modifiers

RegEx – metacharacters

- represent one or multiple characters you want to search for in a string
- Some examples of metacharacters:

<code>^</code>	Matches beginning of line
<code>\$</code>	Matches end of line
<code>.</code>	Matches any single character except newline
<code>[...]</code>	Matches any single character in brackets
<code>[^...]</code>	Matches any single character not in brackets
<code>a b</code>	Matches either a or b

RegEx – metacharacters

^	Matches beginning of line.	(?#...)	Comment.
\$	Matches end of line.	(?= re)	Specifies position using a pattern. Doesn't have a range.
.	Matches any single character except newline. Using m option allows it to match newline as well.	(?! re)	Specifies position using pattern negation. Doesn't have a range.
[...]	Matches any single character in brackets.	(?> re)	Matches independent pattern without backtracking.
[^...]	Matches any single character not in brackets	\w	Matches word characters.
re*	Matches 0 or more occurrences of preceding expression.	\W	Matches non-word characters.
re+	Matches 1 or more occurrence of preceding expression.	\s	Matches whitespace. Equivalent to [\t\n\r\f].
re?	Matches 0 or 1 occurrence of preceding expression.	\S	Matches nonwhitespace.
re{ n}	Matches exactly n number of occurrences of preceding expression.	\d	Matches digits. Equivalent to [0-9].
re{ n,}	Matches n or more occurrences of preceding expression.	\D	Matches nondigits.
re{ n, m}	Matches at least n and at most m occurrences of preceding expression.	\A	Matches beginning of string.
a b	Matches either a or b.	\Z	Matches end of string. If a newline exists, it matches just before newline.
(re)	Groups regular expressions and remembers matched text.	\z	Matches end of string.
(?imx)	Temporarily toggles on i, m, or x options within a regular expression. If in parentheses, only that area is affected.	\G	Matches point where last match finished.
(?-imx)	Temporarily toggles off i, m, or x options within a regular expression. If in parentheses, only that area is affected.	\b	Matches word boundaries when outside brackets. Matches backspace (0x08) when inside brackets.
(?: re)	Groups regular expressions without remembering matched text.	\B	Matches nonword boundaries.
(?imx: re)	Temporarily toggles on i, m, or x options within parentheses.	\n, \t, etc.	Matches newlines, carriage returns, tabs, etc.
(?-imx: re)	Temporarily toggles off i, m, or x options within parentheses.	\1...\9	Matches nth grouped subexpression.
		\10	Matches nth grouped subexpression if it matched already. Otherwise refers to the octal representation of a character code.

RegEx – FLAGS options

- Regex are case sensitive (by default)!
- `re.I` `IGNORECASE`
→ perform case-insensitive matching
- `re.S` `DOTALL`
→ make the `'.'` special character match any character at all, including a newline
Default without this flag= `'.'` will match anything except a newline.
- `re.M` `MULTILINE`
→ the pattern character `'^'` matches at the beginning of the string and at the beginning of each line (immediately following each newline); and the pattern character `'$'` matches at the end of the string and at the end of each line (immediately preceding each newline).
→ Default (without this flag)= `'^'` matches only at the beginning of the string, and `'$'` only at the end of the string and immediately before the newline (if any) at the end of the string.
- `re.X` (`VERBOSE`), `re.L` (`LOCALE`), `re.U` (`UNICODE`)

RegEx – re.split

- Split string by the occurrences of pattern

`re.split(pattern, string, maxsplit=0, flags=0)`

<code>pattern</code>	→ pattern to search
<code>string</code>	→ string to search in
<code>maxsplit</code>	→ maximum number of splits (0 = infinite)
<code>flags</code>	→ optional flags/modifiers

RegEx – re.sub

- `re.sub(pattern, repl, string, count=0, flags=0)`
- Return the string obtained by replacing the leftmost non-overlapping occurrences of pattern in string by the replacement repl.

pattern	→ pattern to search
repl	→ string to substitute the pattern with
string	→ string to search in
count	→ maximum number of substitutions (0 = infinite)
flags	→ optional flags/modifiers

RegEx – re.sub - Groups

- Groups are very handy for substitutions
- Use the following grouping on the line below:

```
re.sub("(gr[ae]y)", "\g<1>blue", line)
```

```
line = "My computer should be grey and my car should also be gray"
```

\g<1> stands for group 1 = the first group between ()

\g<2> stands for group 2, etc..

- Try to understand what happens in the following example:

```
re.sub("(gr[ae]y)(\D*)(gr[ae]y)", "\g<1>blue\g<2>not \g<3>black", line)
```

RegEx – re.search

- re.search (pattern, string, flags=0)

Scan through string looking for the first location (only) where the regular expression pattern produces a match, and return a corresponding **MatchObject instance**.

```
import re
line = "TGCATAGCGAATTCGAGCGT"
match_output = re.search("GAATTC",line)
print match_output

<_sre.SRE_Match object at 0x10c2827e8>

if match_output:
    print "GAATTC site found!"

print match_output.group()      # GAATTC
print match_output.start()      # 8
print match_output.end()        # 14
print match_output.span()       # (8, 14)
```

RegEx – other methods

- `re.subn(pattern, repl, string, count=0, flags=0)`

Return the string obtained by replacing the leftmost non-overlapping occurrences of `pattern` in `string` by the replacement `repl`. If the `pattern` isn't found, `string` is returned unchanged. Returns a tuple with the number of changes.

- `re.escape(string)`

Return string with all non-alphanumerics backslashed

- Returns Objects

- `re.match(pattern, string, flags=0)`

If zero or more characters at the beginning of `string` match the regular expression `pattern`, return a corresponding `MatchObject` instance

- `re.finditer(pattern, string, flags=0)`

Return an iterator yielding `MatchObject` instances over all non-overlapping matches for the RE `pattern` in `string`

- Compile regex (generally for speed purposes)

- `re.compile(pattern, flags=0)`

Compile a regular expression `pattern` into a regular expression object

Regular Expression



Anchors		Sample Patterns	
^	Start of line +	([A-Za-z0-9-]+)	Letters, numbers and hyphens
\A	Start of string +	(\d{1,2}\V\d{1,2}\V\d{4})	Date (e.g. 21/3/2006)
\$	End of line +	([^\s]+(?:\.(jpg gif png))\.\2)	jpg, gif or png image
\Z	End of string +	(^[1-9]{1}\$ ^[1-4]{1}[0-9]{1}\$ ^50\$)	Any number from 1 to 50 inclusive
\b	Word boundary +	(#?([A-Fa-f0-9]){3}((([A-Fa-f0-9]){3})?))	Valid hexadecimal colour code
\B	Not word boundary +	((?=\.*\d)(?=\.*[a-z])(?=\.*[A-Z]).{8,15})	8 to 15 character string with at least one upper case letter, one lower case letter, and one digit (useful for passwords).
\<	Start of word	(\w+@[a-zA-Z_]+?\.[a-zA-Z]{2,6})	Email addresses
\>	End of word	(\<(/?[^\>]+)\>)	HTML Tags
Character Classes			
\c	Control character	Note <i>These patterns are intended for reference purposes and have not been extensively tested. Please use with caution and test thoroughly before use.</i>	
\s	White space		
\S	Not white space		
\d	Digit		
\D	Not digit		
\w	Word		
\W	Not word		
\xhh	Hexadecimal character hh		
\Oxxx	Octal character xxx		
POSIX Character Classes			
[[:upper:]]	Upper case letters	Quantifiers	
[[:lower:]]	Lower case letters		
[[:alpha:]]	All letters		
[[:alnum:]]	Digits and letters		
[[:digit:]]	Digits		
[[:xdigit:]]	Hexadecimal digits		
[[:punct:]]	Punctuation		
[[:blank:]]	Space and tab		
[[:space:]]	Blank characters		
[[:cntrl:]]	Control characters		
[[:graph:]]	Printed characters	Ranges	
[[:print:]]	Printed characters and spaces		
[[:word:]]	Digits, letters and underscore		
[[:upper:]]	Upper case letters		
[[:lower:]]	Lower case letters		
[[:alpha:]]	All letters		
[[:alnum:]]	Digits and letters		
[[:digit:]]	Digits		
[[:xdigit:]]	Hexadecimal digits		
[[:punct:]]	Punctuation		
[[:blank:]]	Space and tab		
[[:space:]]	Blank characters		
[[:cntrl:]]	Control characters		
[[:graph:]]	Printed characters		
[[:print:]]	Printed characters and spaces		
[[:word:]]	Digits, letters and underscore		
Assertions		Special Characters	
?=	Lookahead assertion +		
?!	Negative lookahead +		
?<=	Lookbehind assertion +		
?!= or ?<!	Negative lookbehind +		
?>	Once-only Subexpression		
?()	Condition [if then]		
?()	Condition [if then else]		
?#	Comment		
Note		Items marked + should work in most regular expression implementations.	
		String Replacement (Backreferences)	
		Metacharacters (must be escaped)	
		Available free from AddedBytes.com	

Exercises