

# Regular expressions – Cheat-sheet and examples

(V. Nicosia – 2021)

## Definitions

**Character and string.** Characters are the smallest components of a string, i.e., single letters, digit, punctuation symbols, etc. A string is a sequence of characters. The string “Queen!” consists of the six characters 'Q', 'u', 'e', 'e', 'n', '!'.

**Regular expression.** A regular expression (in short **regex**) is a pattern that describes a set of strings. We say that a string “matches” a given regular expression if it adheres to the pattern defined by the regex.

**Meta-character.** Some characters have a special meaning when used in regular expressions, i.e., they do not directly represent the corresponding character but are used to signify something else. Meta-characters include: '.' (dot), '\*', '+', '\\', '(', ')', '[', ']', '{', '}', '?', '^', '\$'. If you want to indicate a meta-character as literal (for instance because you are looking for a string that contains one of them), you need to precede it by a '\\'. For instance, the pattern “a\\\*b” indicates a string formed by the three characters 'a', '\*', 'b', while the pattern “a\*b” is a regular expression that matches a sequence of zero or more 'a' followed by a 'b' (no '\*' in there...).

## Patterns

Groups and ranges				
	Meaning	grep -E	grep -P	Python re
.	any character	Y	Y	Y
[abcd]	one character among those specified within []	Y	Y	Y
[^ab]	none of the characters specified within []	Y	Y	Y
[b-g]	any of the characters in a certain range (here, the characters between 'b' and 'g', both included)	Y	Y	Y
(bar)	a group of characters (here the sequence “bar”)	Y	Y	Y

Quantifiers				
	Meaning	grep -E	grep -P	Python re
*	zero or more of whatever preceeds it	Y	Y	Y
+	one or more of whatever preceeds it	Y	Y	Y
?	zero or one of whatever preceeds it	Y	Y	Y
{2}	exactly 2 of whatever preceeds it	Y	Y	Y
{2,6}	between 2 and 6 times of whatever preceeds it	Y	Y	Y
{4,}	four or more of whatever preceeds it	Y	Y	Y

Anchors				
	Meaning	grep -E	grep -P	Python re
<code>^</code>	the beginning of string or line	Y	Y	Y
<code>\$</code>	the end of string or line	Y	Y	Y
<code>\A</code>	start of string	N	Y	Y
<code>\Z</code>	end of string	N	Y	Y
<code>\b</code>	word boundary	N	Y	Y
<code>\B</code>	not a word boundary	N	Y	Y
<code>&lt;</code>	start of word	N	Y	Y
<code>&gt;</code>	end of word	N	Y	Y

Character classes				
	Meaning	grep -E	grep -P	Python re
<code>\c</code>	a control character	N	Y	Y
<code>\s</code>	a white space (including tabs)	N	Y	Y
<code>\S</code>	not a white space	N	Y	Y
<code>\d</code>	a digit	N	Y	Y
<code>\D</code>	not a digit	N	Y	Y
<code>\w</code>	a word	N	Y	Y
<code>\W</code>	not a word	N	Y	Y
<code>\x</code>	hexadecimal digit	N	Y	Y
<code>\O</code>	octal digit	N	Y	Y
<code>[:digit:]</code>	a digit	Y	Y	N
<code>[:alpha:]</code>	a letter	Y	Y	N
<code>[:alnum:]</code>	an alphanumeric character (letter or digit)	Y	Y	N
<code>[:space:]</code>	a blank space	Y	Y	N
<code>[:cntrl:]</code>	a control character	Y	Y	N
<code>[:lower:]</code>	a lowercase letter	Y	Y	N
<code>[:upper:]</code>	an uppercase letter	Y	Y	N
<code>[:punct:]</code>	a punctuation symbol	Y	Y	N

## Reading out regular expressions

You find below some relevant examples of regular expressions, together with the way each of them would be expressed in plain English. The same list of regular expressions is reported in the table on the next page, with examples of which strings they do or do not match.

Expression	reads out as...
<code>arch</code>	the string “arch”, i.e., an ‘a’, followed by an ‘r’, followed by a ‘c’, followed by an ‘h’
<code>...</code>	a sequence of any three characters
<code>ab.c</code>	the string “ab”, followed by any single character, followed by the character ‘c’
<code>^arch</code>	a string that starts with the substring “arch”
<code>ber\$</code>	a string that ends with the substring “ber”
<code>^co.*ler\$</code>	a string that starts with “co”, followed by zero or more characters, and ends with “ler”
<code>a.*t</code>	the character ‘a’, followed by zero or more occurrences of any single character, followed by a ‘t’
<code>(at){2,}</code>	the string “at” repeated two or more times
<code>at.+(at){1,}</code>	a string that contains the string “at”, followed by at least one of any character, followed by the string “at” repeated one or more times
<code>^[aeiou]</code>	a string that starts with a lowercase vowel
<code>[aeiou]{2,}</code>	a string that contains a sequence of two or more vowels
<code>^[^aeiouAEIOU]</code>	a string that does not start with a vowel (either lowercase or uppercase). Equivalently, a string that starts with a consonant! (not a vowel)
<code>[0-9]+.*[g-j]{2}\$</code>	A string that contains a sequence of one or more digits (a number), then zero or more characters, and ends with two contiguous letters from the set [‘g’, ‘h’, ‘i’, ‘j’]
<code>(jo ra)b</code>	a string that contains either the string “jo” or the string “ra”, each followed by the letter ‘b’
<code>1[0-9]{2}</code>	a string containing a ‘1’ followed by any two digits
<code>^[0-9]+\$</code>	a string that consists exclusively of digits (at least one!)
<code>^([-a-zA-Z0-9_]+\.){2}uk\$</code>	a string consisting of exactly two sequences of the form “one or more letters or digits or ‘-’ (dash) or ‘_’ (underscore) followed by a “.” (dot)”, then ending with “uk”

## Examples

The same set of regular expressions reported in the table above, with relevant examples of which strings they would or would not match.

Expression	Matches	Does not match
arch	" <u>March</u> ", " <u>archery</u> ", " <u>starch</u> "	"ARCH", "arc", "Arch"
...	" <u>matter</u> ", " <u>abc</u> ", " <u>a b</u> ", " <u>b2c</u> "	"ab", "" (empty string), "13"
ab.c	" <u>abacus</u> ", " <u>slab chain</u> ", " <u>abscess</u> ", " <u>Arabic</u> "	"abcd", "chamber", "Abacus"
^arch	" <u>arch</u> ", " <u>archery</u> ", " <u>archeology</u> "	"March", "under the arch", "Archery"
ber\$	" <u>amber</u> ", " <u>September</u> ", " <u>number</u> ", " <u>chamber</u> "	"beer", "bergamot", "numbered", "cumbersome"
^co.*ler\$	" <u>compiler</u> ", " <u>cooler</u> ", " <u>cobbler</u> "	"toddler", "cooled"
a.*t	" <u>at</u> ", " <u>hat</u> ", " <u>yachts</u> ", " <u>compartment</u> ", " <u>aesthetic</u> ", " <u>breakfast</u> ", " <u>alphabet</u> "	"bit", "mash", "At", "haT"
(at){2,}	" <u>datatype</u> ", " <u>catatonic</u> "	"latter", "athletics", "wombat"
at.+(at){1,}	" <u>bathmat</u> ", " <u>attenuated</u> ", " <u>saturation</u> "	"datasets", "datatype"
^[aeiou]	" <u>aerospace</u> ", " <u>i</u> ", " <u>union</u> "	"Ihaho", "trampoline", "22catch"
[aeiou]{2,}	" <u>armour</u> ", " <u>Lilliputian</u> ", " <u>raining</u> ", " <u>queue</u> "	"TRAIN", "last"
^[^aeiouAEIOU]	" <u>train</u> ", " <u>random</u> ", " <u>Hotel</u> "	"indigo", "orange", "Apples"
[0-9]+.*[g-j]2\$	"write <u>22 emoji</u> ", " <u>120-laugh</u> ", "give me <u>1 egg</u> "	"have a laugh", "I got 3 pigs", "You need to go through it"
(jo ra)b	"my new <u>job</u> is fun", " <u>crabs</u> ", "I found a <u>scarab</u> "	"jab", "rub", "jorab"
1[0-9]{2}	" <u>178</u> ", " <u>alpha-102-bravo</u> ", " <u>197865</u> ", " <u>231462</u> ", "Capt. Sir Tom Moore died at age <u>100</u> "	"22.87", "1st position", "232462"
^[0-9]+\$	" <u>7</u> ", " <u>18764562340</u> ", " <u>888</u> "	"alpha33", "3 and me", "78 times 22"
^([-a-zA-Z0-9_]+\.){2}uk\$	" <u>yahoo.co.uk</u> ", " <u>your-email.co.uk</u> ", " <u>data.gov.uk</u> ", " <u>02_real.mean.uk</u> "	"www.yahoo.co.uk", "mail.uk.gov", "lab.inria.fr", "sunsite.unc.edu"