

NLP Sheet 1

1) Write a regular expression that matches any element in the set of

⑥

a) All Alphabetic Strings

→ Means any sequence of alphabetic characters

→ Clearly

$/[a-zA-Z]^+ /$

} one or more
alphabetic
characters

b) All lowercase alphabetic strings ending in b

→

$/[a-z]^* b /$

a sequence of
alphabetic char.
(can be nothing)

ends with
b

Notice that given
abcb → The * matches the
longest string it can

c) All strings from the alphabet {a, b} such that each a is immediately preceded by and followed by b

→ This means that such sequence is only allowed

① to start with b (bbb... is a general beginning)

one or more

→ If an a ever wants to appear then what

② must come after is any ≥ 1 no. of bs

(bbbba bbb...)

This scenario can repeat anytime later (e.g. →)

③ for any no. of times

Combining ①, ②, ③

$$/ b + (ab^+)^* /$$

① ② ③
applies on ()

d) All binary strings with at least 4 1s

- Should be a sequence of 0s and 1s
- 1 should appear 4 times or more at least

* Consider the 1st Four 1s in a random binary string. They must be separated by 0s from each other. ① ②
any no. of

* The rest of the string can be any binary string ②

$$/ 0^* 1 0^* 1 0^* 1 0^* 1 [01]^* /$$

① ②

This is clearly equivalent to the expression

$$/ \underbrace{(0^* 1)^4}_{\text{①}} \underbrace{[01]^*}_{\text{②}} /$$

* Another way to think about is that if we spot 4 1s anywhere in a binary string then what separates them is an arbitrary binary string (of any length)

$/ ([01]^* 1) \{4\} [01]^* /$

--1--1--1--1--
↓ ↓ ↓ ↓ ↓
any binary
string or nothing

* Another way to think about it is that all the ones in a random binary string are pair-wise separated by any no. of zeros.

i.e. $/ (0^* 1) \{4\} 0^* /$

--1--1--1--1--1--
↑ ↑ ↑ ↑ ↑
all thots
↓ ↓ ↓ ↓ ↓
zeros
(any no.)

* which can be also written as

 $/ (0^* 1 0^*) \{4\} /$

. we used 1+ in tut.

" In a random binary string, every single 1 you find will be separated from the next by zeros. In our case, this must happen 4 times (any no.) (or more)

- This hopefully covers the answers we considered in the tutorial

e) All binary strings where the no. of zeros is a multiple of 3.

→ Need 0, 3, 6, 9, ... zeros

→ They are separated from each other by 1s

Don't want to "match"
 $/ (1^* 0 1^* 0 1^* 0 1^*)^+ /$

- This assumes no zeros isn't allowed

• We can write the previous as $/((1^* 0 1^*)^3)^+ /$

→ To allow that there can be no zeros in the expression (just 1s)

$/((1^* 0 1^*)^3)^+ | 1 + /$

→ Notice that we never resorted to boundaries because we were never constrained to do so (boundaries help guarantee the matched words are surrounded by the right characters)

2) Define a word as

→ Sequence of alphabetic characters

→ Separated from other words through white space, any relevant Punct., line break's, ...

Does
`\\b`
still
work as
a
separator
here?

• In regex a word boundary is a position where a non-regex word occurs

↳ [a-zA-Z0-9_]

(anything that isn't a letter, no. or _ can be there)

↳ includes white space, relevant Punct., and more that wasn't strictly disallowed

Thus, also the word's definition has changed,
`\\b` still works as a separator.

* Write Regex to find elements belonging to the set of

a) All Strings with two consecutive repeated words in the same case

$/ \backslash b ([a-zA-Z]^+) \backslash S+ (\backslash 1) \backslash b /$

This also assumes that the two words can be only separated by spaces (and not tabs).
works.

Capturing group

{ • assume any amount of space between them is okay (examples show -1)

Perhaps, $[w]^+$ instead of $\backslash S^+$ can work for that

b) All Strings that start the beginning of the line with with an integer and that end the end of the line with a word.

$/ ^ \backslash d+ \cdot * [a-zA-Z]^+ \$ /$

Start with $\backslash d^+$ ends with a
an int of 1 or more digits word

but Notice

→ It's sufficient to check that it starts with one integer (rest in $\cdot *$)

→ Integers can be -ve } add (-) ?

→ to be a word ① must be preceded by a word boundary but $\cdot *$ allows anything (even no.5)

Hence, need to rewrite as

$/ ^ (-)? \d \cdot * \b [a-zA-Z] + \$ /$

↓
condrop
()

→ last char. in $\cdot *$ must be a separator.

c) All Strings that have both the word grotto and the word raven in them (neither should be a SubString of another word as you'd expect)

$/ \cdot * \b grotto \b \cdot * \b raven \b \cdot * /$

- This clearly only catches them if grotto appears before raven
 - Need to 'OR' the opposite order.
(Pipe)

<bonus>

→ Can we generalize. What if we're checking for existence of more words (too many permutations to OR)



$^ \cdot * (?=\b one\b) \cdot * (?=\b two\b) \cdot * (?=\b three\b) \cdot *$

• existence of 'one', 'two', 'three' checked regardless to Order

• Thanks Saad <3

3) Write a Single Regex that matches all of

Colours } optional u ①
Colors
They're Colours } optional they're ②
they're Colors
they are Colours } instead of they're, they are ③
they are Colours } is also accepted

① /Colo(u)?rs/

② / (they're)?Colo(u)?rs/

③ / (they're)?Colo(u)?rs/

* in the tutorial we used \s for the two spaces.

Thanks 23