# Regex

* match/check Series of characters (email validation: .Com, @, URL)

## Chan Set

[A-z a-z] : A-Z, a-z

[^a-z] : except a-z

[0-9] : 0-9

## Repeating:

[a-z]+ : Atleast one (1/more)

[a-z]{10} : 10 times exactly

[a-z]{5,8} : B/w 5 to 8

[a-z]{5,} : 5 to ∞

## Metacharacters:

\d : digits [0-9]

\w : [a-z][A-z][0-9] -

\S : Any whitespaces except newline (Spaces, tabs)

\t : Tabs only

\ : Escape character

\d{3} \S \w{5} : 3 digits, 1 space, 5 letters

## Special characters:

+ : 1/more

\ : Escape character

[] : chan set

[^] : Negate Symbols Inside chan set

? : zero | one

. : Any char (except newline)

* : 0/more

() : group regex

### Start with:

^[a-z]{5}

### Start & end with:   ^[a-z]{5}$

## OR:

p| tyre : [P]tyre : matches 'p' only

(p|t) yre : matches p|t then yre.

eg:

(name) @ (domain). (extension)(.again)

boss @ ninja. com[.uk]
      └→ optional

$\wedge ([a-zA-z0-9\backslash,-]+)@([a-zA-zA-z\backslash d\,-]+)\backslash.([a-zA-zA-z]\{2,8\})$

$(\backslash.[a-zA-z A-z]\{2,8\})?\$$

---

\A(__) : Begin with given char(s)

\b(—) : begin/end with given character(s)

\B(—) : Shouldn't begin/end ''

\d   : [0-9]

\D   : [^ 0-9]

\s   : whitespaces

\S   : Non whitespaces

\w   : [a-zA-z0-9 _]

\W  : Any non-alpha numeric char

\Z  : String ends with given regex

Note: Inside sets

   [+] : No special meaning

    [+, *, ., ?, ', |, (), $, {3}]

## Python functions:

re.findAll() : return non-overlapping matches; order maintained

re.Compile() : Compile to pattern (Search, Substitute)

re.Split() : Upon finding Split

re.Sub() : Replace Content

re.Subn() : Tuple with Count of replacement, new String (replaced)

re.escape() : Return String : remove all non-alphanumerics backslashed (match String have regex metachars)

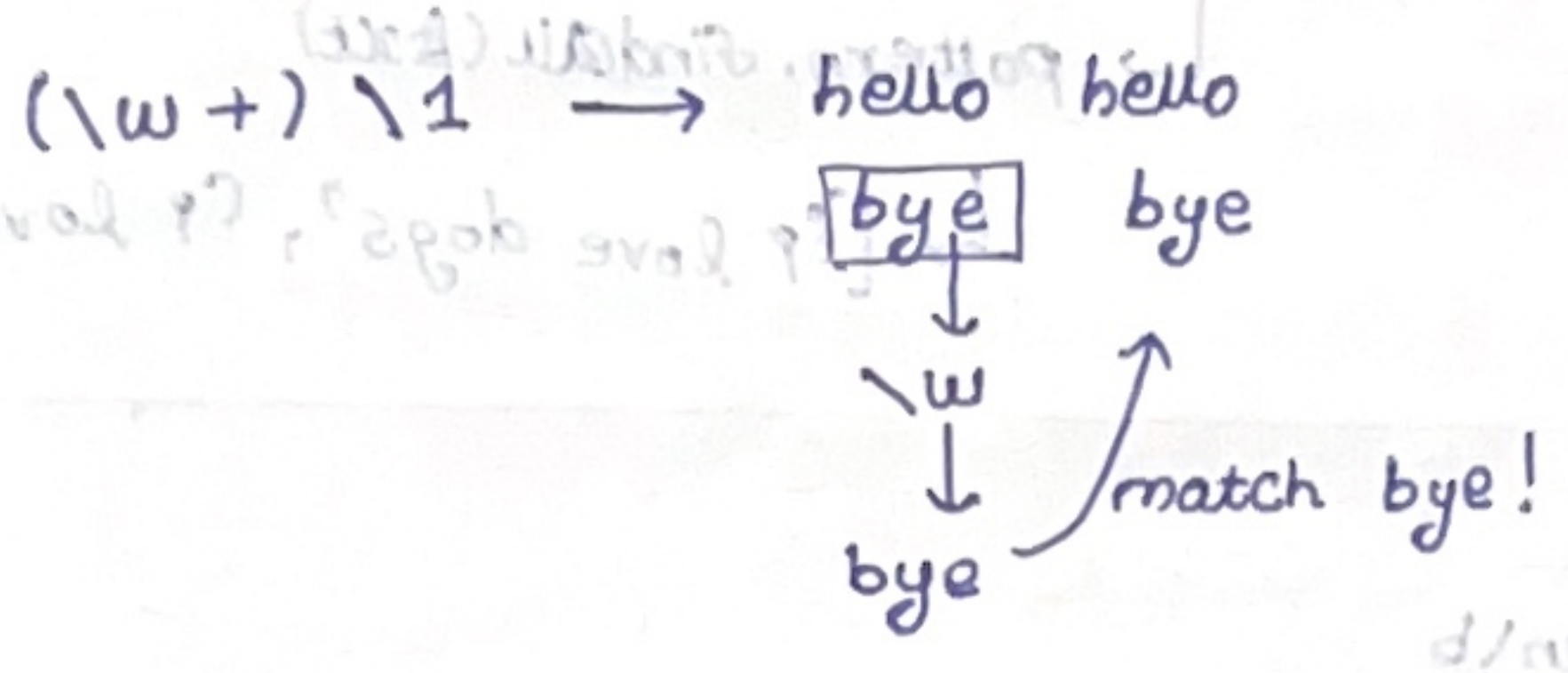re.Search() : Return None / re (more than extracting : used for Testing)!

# J' Sticks (Start Searching from 'lastIndex' property)

## Flags:

- ?: Ignore Case
- M: multiline
- S: make . match newline also
- U: maker $\{$ \w \W \b \B $\}$ follow unicode rules
- L:        "           Locale (lang, time, ..) : Not recomm!
- A:        "           ASCII only
- X: Allow Comment in regex
- DEBUG: Get info about regex Compilation pattern!
- g: global

## Back Referencing:

* Match String using o/p of regex group 1

$$(\w+) \ \backslash 1 \longrightarrow \quad hello \ hello$$

$\boxed{bye}$  bye

$\downarrow$ \w

$\downarrow$  /match bye!

bye

## Named groups:

$(?P<firstName> \w+) (?P<last> \w+ )$

$\boxed{match.group(`firstName')}$

## Swap using regex:

$$Pattern.sub(``\backslash 2 \ \backslash 1", \ txt)$$

$$Pattern.sub(r``\g<last> \g<first>", txt)$$

## Same First & last name:

`John John'

Pattern = re.Compile (`(\w+) \1")

(or) [ Pattern = re.Compile (`(?P<first> \w+) (?P=first)") ]

$\boxed{Pattern.match(txt)}$

match.groupdict() $\longrightarrow$ {`first' : `John', `last' : `John'}

# Non Capturing groups:

      ★ Eg: alternation (Something can't be altered)

'I love (Cats|dogs)' ⟶ ~~pattern~~. Find All(txt)

                    ['cats', 'dogs'] ⟶ returns group matches

## case:

     I want entire String instead of groups : Not matching groups

     ┌─────────────────────────────────────────────────────┐
     │ (?: pattern) ⟶ Non matching group! │
     └─────────────────────────────────────────────────────┘

           re. Compile ('I love (?:cats/dogs)')

                  ┗➤ pattern. findAll (txt)

                     ┗➤ ['I love dogs', 'I love cats']

---

# word boundaries (\b):

         \bgreen\b

      'red green blue' ⟶ ['green']

## Zero-width assertions: characters indicating positions ($\wedge$, \$) rather than content (like assertion)

## Powerful Zero-width assertion: Look around
                                Look ahead
                                Look behind

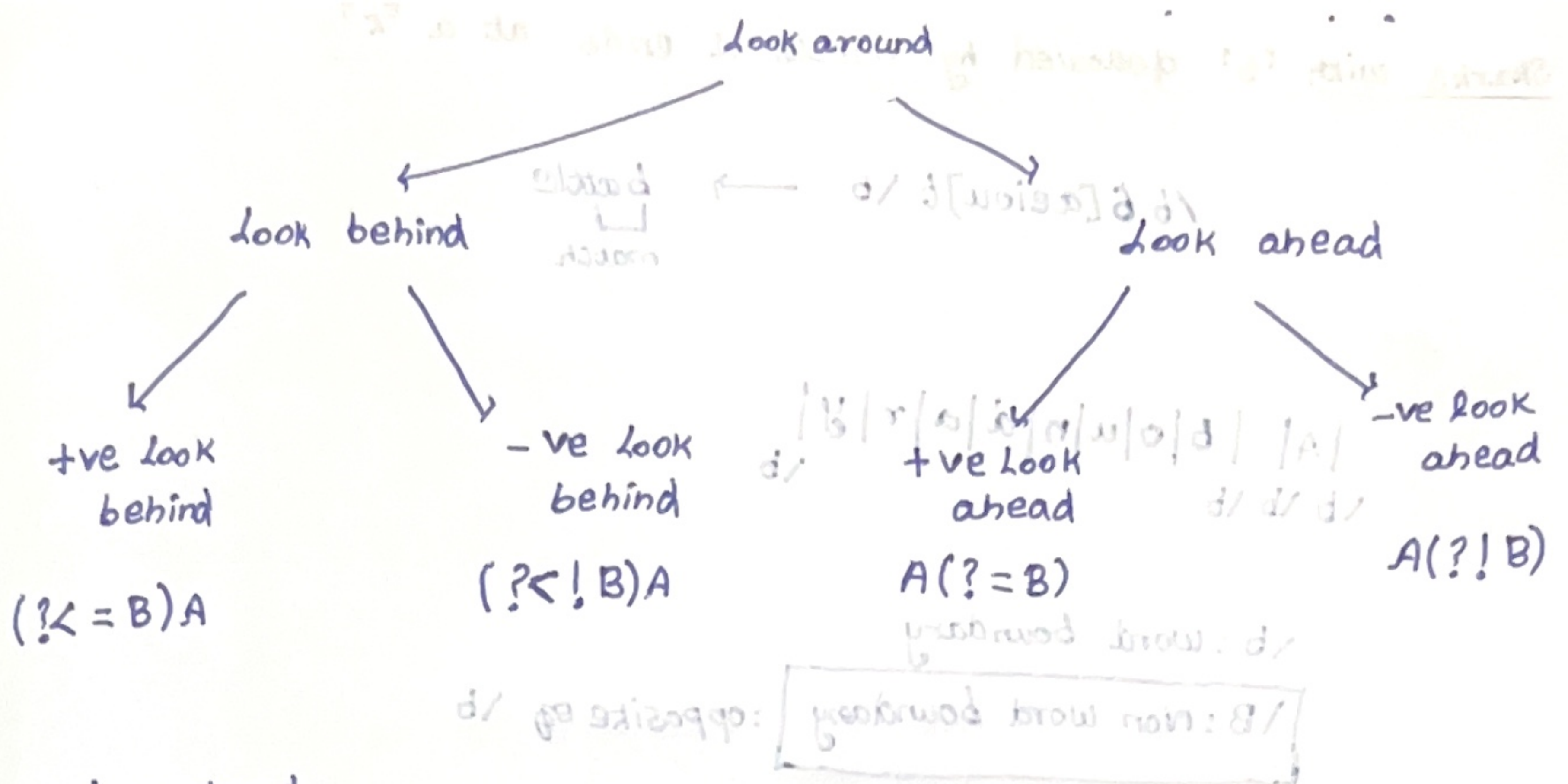## Look around: Certain regex before & after a word/regex!

                  Look around (2 parts)

                  ↙               ↘

        Actual             Non-consuming
        expression        expression
        (final result)     (Just to check position
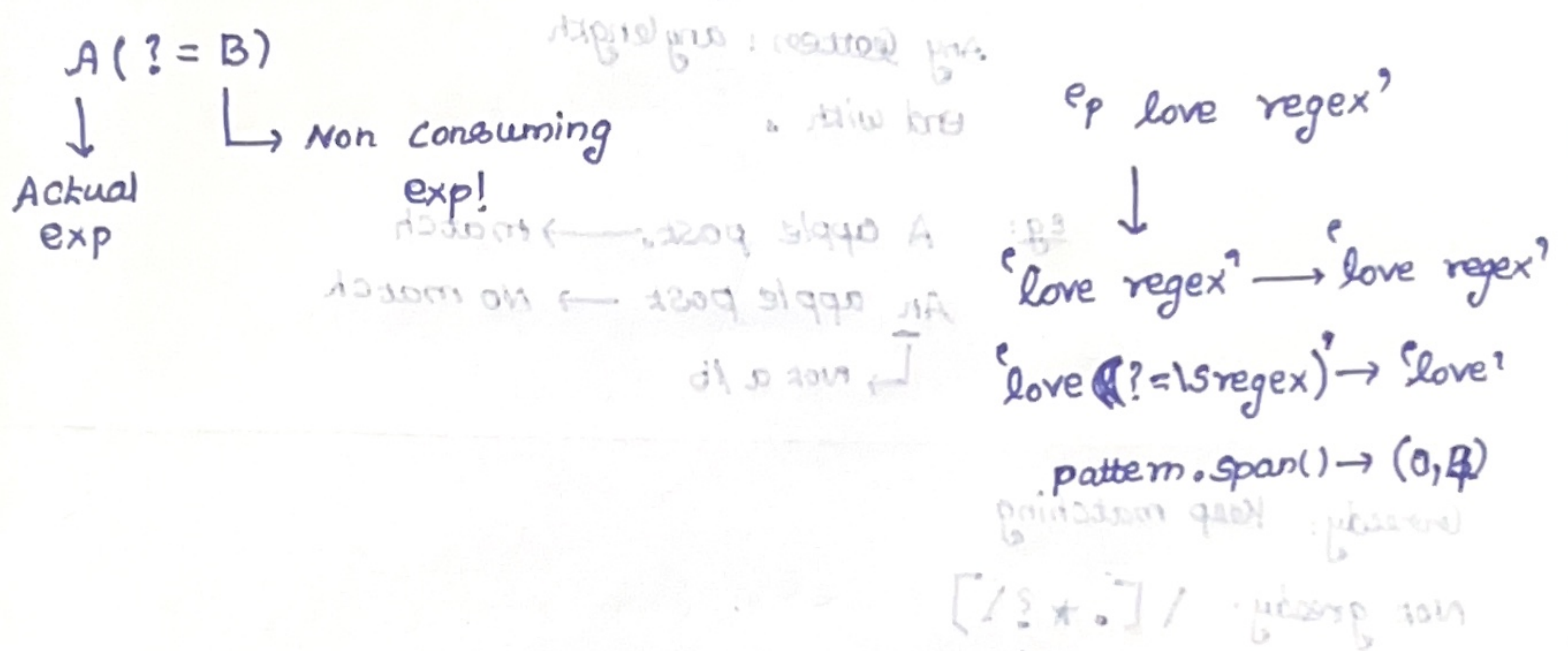                          like can succeed: start/
                               end)

                              ↓

                        doesn't match: move to
                         next char & do branching
                           matching again.

look around

look behind ← → look ahead

look behind:
- +ve look behind
- −ve look behind

look ahead:
- +ve look ahead
- −ve look ahead

+ve look behind
$(?<=B)A$

−ve look behind
$(?<!B)A$

+ve look ahead
$A(?=B)$

−ve look ahead
$A(?!B)$

## +ve Look ahead:

* check for pattern ahead of a regex expression

$A(?=B)$
↓   ↳ Non consuming exp!
Actual exp

'p love regex'
↓
'love regex' → 'love regex'
'love(?=\sregex)' → 'love'

pattern.span() → (0,4)

eg:
'my favourite colors are red, green, and blue?'

All words followed by comma (or) fullstop:

$\w+(?=[,.])$    (_____)
$\w+(?=,|\.)$    (_____)

## Negative look ahead:
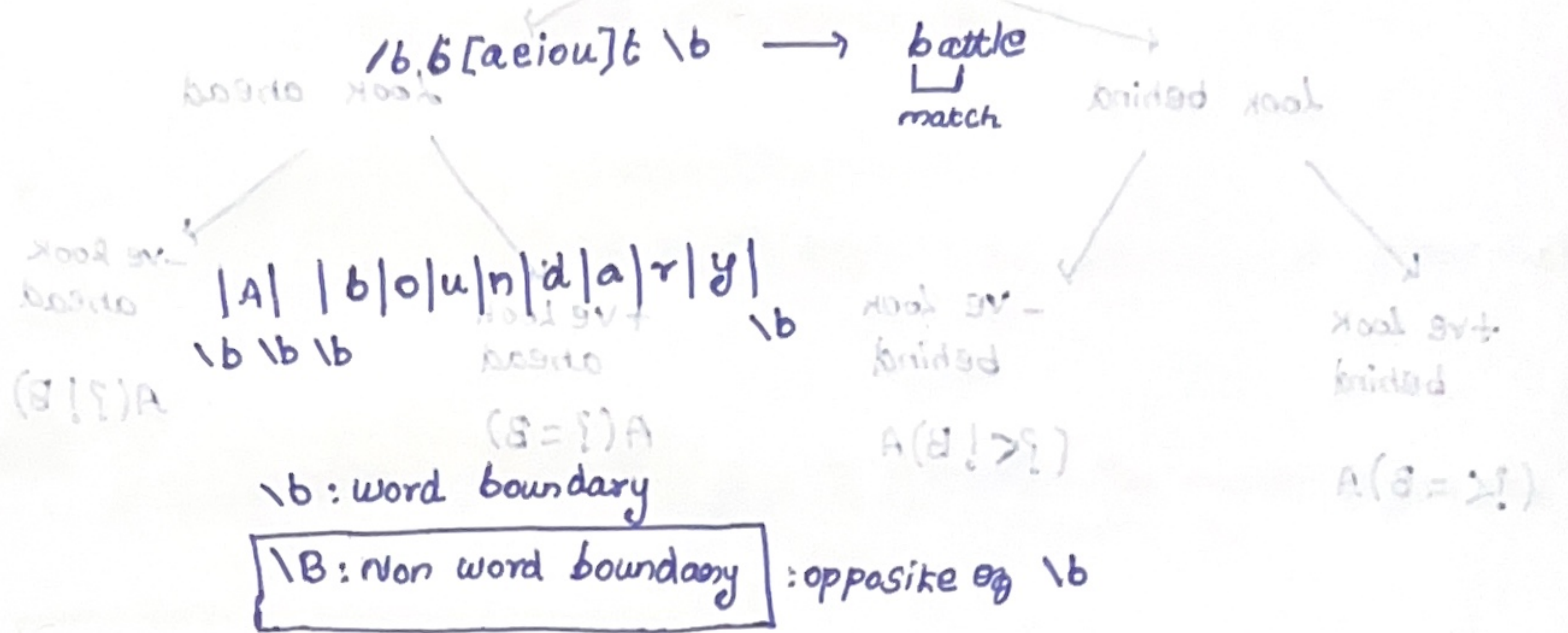
* match only not followed by a regex
eg: 'love' not followed by 'regex'

'p. love physics, p love regex'

'love(?!\s regex)' ⟶ span? (2,6)

| 'Same with look behind' |

Starts with 'b' followed by vowels & ends at a 'k'

/b.b[aeiou]t \b ⟶ battle
⎵
match

|A| |b|o|u|n|d|a|r|y|
\b \b \b                    \b

\b : word boundary

\B : Non word boundary : opposite of \b

^A\b.*\.$ ⟶ Starts with A
word boundary
Any pattern : any length
end with .

eg: A apple post. ⟶ match
An apple post ⟶ No match
↳ Not a \b

Greedy: Keep matching

Not greedy: \[.*?\]
↳ Not greedy

greedy: match until ']' no longer available.

[Google] (____ure____) [test]
[?bm] (____ure____)

\[.*\] ⟶ '[Google](ure)[test]'

\[.*$\] ⟶ '[Google]', '[test]'

Group within groups:

(cat (\d\d)+)+ ⟶ Cat cat30cat5084
Cat 54 Cat 30

Note: Subgroups : tedious task (make sure of detail
such as parenthesis, a minute
change : different result.)

'The Code Us 45cat7814Cat' ⟶ 45cat7814 cat
e
    " 45cat_7814CatCat' ⟶ match nothing


?: and ? = (confusion)
_____

        ★ whenever ?= (equal) we are talking about assertion.