## Agenda

- Regex
- Simple pattern Matching
- Meta Characters
- Anchors
- Quantifiers and Groups
- Python implementation of Regex with re library

## Regex

Regex can help in identifying complex Patterns in Documents

- Data Cleaning
- Data Validation
- Masking Information

858683 6627 ⇒ ****** ** 6627

1) Case - sensitive     `abc` != `ABC`
2) Order matters     `abc` != `cba`

`.`

3) Special character that matches Everything Except newline.

\|     3) Escape character / backslash

`\.`  4) metachar

# Meta Characters

\d : Matches all digits

\D : | Matches all digits

\w : Matches Alphanumeric and _

\W : ! \w

\s : matches Whitespace chars

\S : ! \s

```
.           - Any Character Except New Line
\d          - Digit (0-9)
\D          - Not a Digit (0-9)
\w          - Word Character (a-z, A-Z, 0-9, _)
\W          - Not a Word Character
\s          - Whitespace (space, tab, newline)
\S          - Not Whitespace (space, tab, newline)
```
meta chars

```
\b          - Word Boundary
\B          - Not a Word Boundary
^           - Beginning of a String
$           - End of a String
```
Anchors

```
[]          - Matches Characters in brackets
[^ ]        - Matches Characters NOT in brackets
|           - Either Or
( )         - Group
```

# Anchors

1. Anchors Don't match any characters

2. They match invisible positions before or after characters

① \b : boundary of a word
   - Prefix

② \B : suffix

③ ^ → matches Beginning of String

④ $ → Matches end of String

Pattern$

# Character Set

[. -]  ⊘ matches either .  or -

[19]x    [1 2 3 4 5 6 7 8 9] ✓

↓

[1 - 9]

Ex

[a b]

[. —]

* numeric character set

⊘ [1 2 3 4 5 6]

⊙ [1 - 6]    ⊙ [start — end]

⊘ [1 - 9]

* Alphabet charset

[a - z]    [A - Z]    [a - z A - Z]

* Negation in character-set

[^a-z] ⇒ match everything but lowercase char

* Note: inside a character Set
   ∧ works as negation of
   character Set

---

## Quantifiers

⇒ match specific repetion of characters

① * → 0 or more

② + → 1 or more

③ ? → 0 or One

④ {3} ⇒ Exactly 3

   {3, 4} ⇒ {Min 3, max-4}

$\{3,3\} \rightarrow 3$ or more

|d|d|d| ①   |d $\{3\ 3\}$

Mr|o? → Mr / Mr.

$$\boxed{Group}$$

$$(\qquad)$$

$$\boxed{M} \quad (r \mid s \mid rs) \quad \backslash . \ ? \quad [a-zA-Z]*$$

$$r \quad \checkmark$$

or

$$s \quad \checkmark$$

or

$$rs \quad \checkmark$$

$$(Mr \mid Ms \mid Mrs)$$