# Introduction to regular expressions

## REGULAR EXPRESSIONS IN PYTHON

**Maria Eugenia Inzaugarat**
Data Scientist

# What is a regular expression?

**REG**ular **EX**pression or regex:

*String containing a combination of normal characters and special metacharacters that describes patterns to find text or positions within a text*

r'st\d\s\w{3,10}'

2. What is a regular expression?
A regular expression, or regex, is a string that contains a combination of normal and special characters that describes patterns to find text within a text.
In Python, the r at the beginning indicates a raw string. It is always advisable to use it.

# What is a regular expression?

REGular EXpression or regex:

*String containing a combination of **normal characters** and special metacharacters that describes patterns to find text or positions within a text*

$$r'st\backslash d\backslash s\backslash w\{3,10\}'$$

- Normal characters match themselves ( `st` )

3. What is a regular expression?
We said that a regex contains normal characters, or literal characters we already know. The normal characters match themselves. In the case shown in our slide, st exactly matches an s followed by a t.

# What is a regular expression?

REGular EXpression or regex:

*String containing a combination of normal characters and **special metacharacters** that describes patterns to find text or positions within a text*

$$r'st\backslash d\backslash s\backslash w\{3,10\}'$$

- Metacharacters represent types of characters ( `\d` , `\s` , `\w` ) or ideas ( `{3,10}` )

4. What is a regular expression?
They also contain special characters. Metacharacters represent types of characters. Let's look one by one as they appear in the slide. Backslash d represents a digit.

# What is a regular expression?

> REGular EXpression or regex:

> *String containing a combination of normal characters and **special metacharacters** that describes patterns to find text or positions within a text*

$$r'st\backslash d\backslash s\backslash w\{3,10\}'$$

- Metacharacters represent types of characters ( `\d` , `\s` , `\w` ) or ideas ( `{3,10}` )

5. What is a regular expression?
backslash s a whitespace,
6. What is a regular expression?
backslash w a word character. They also represent ideas, such as location or quantity.

7. What is a regular expression?
In the example, 3 and 10 inside curly braces indicates that the character immediately to the left, in this case backslash w, should appear between 3 and 10 times.

# What is a regular expression?

**REGular EXpression or regex:**

*String containing a combination of normal characters and **special metacharacters** that describes patterns to find text or positions within a text*

r'st\d\s\w{3,10}'

- Metacharacters represent types of characters ( `\d` , `\s` , `\w` ) or ideas ( `{3,10}` )

# What is a regular expression?

**REGular EXpression or regex:**

*String containing a combination of normal characters and **special metacharacters** that describes patterns to find text or positions within a text*

$$r'st\backslash d\backslash s\backslash w\{3,10\}'$$

- Metacharacters represent types of characters ( `\d` , `\s` , `\w` ) or ideas ( `{3,10}` )

9. What is a regular expression?
As a data scientist, you will use pattern matching to find and replace specific text. To validate strings such as passwords or email addresses. Why use regex? They are very powerful and fast. They allow you to search complex patterns that would be very difficult to find otherwise.

# What is a regular expression?

**REGular EXpression or regex:**

*String containing a combination of normal characters and special metacharacters that describes **patterns** to find text or positions within a text*

- Pattern: a sequence of characters that maps to words or punctuation

10. The re module
Python has a useful library, the re module, to handle regex. You can import it as shown in the code. Let's see how it works. To find all matches of a pattern, we use the dot findall method. It takes two arguments: the regex and the string. In the code, we want to find all the matches of hashtag movies in the specified string. The method returns a list with the two matches found.

# What is a regular expression?

**REGular EXpression or regex:**

*String containing a combination of normal characters and special metacharacters that describes patterns* **to find text or positions within a text**

- Pattern matching usage:
  - Find and replace text
  - Validate strings

- Very powerful and fast

# The re module

```
import re
```

- Find all matches of a pattern:



```
re.findall(r"#movies", "Love #movies! I had fun yesterday going to the #movies")
```

```
['#movies', '#movies']
```

10. The re module
Python has a useful library, the re module, to handle regex. You can import it as shown in the code. Let's see how it works. To find all matches of a pattern, we use the dot findall method. It takes two arguments: the regex and the string. In the code, we want to find all the matches of hashtag movies in the specified string. The method returns a list with the two matches found.

# The re module

```
import re
```

- <mark>Split</mark> string at each match:

$$\text{re.split}(r''\text{regex}'', \text{string})$$

```
re.split(r"!", "Nice Place to eat! I'll come back! Excellent meat!")
```

```
['Nice Place to eat', " I'll come back", ' Excellent meat', '']
```

11. The re module
To split a string at each pattern match, we could use the method dot split In the example, we want to split the specified string at every exclamation mark match. It returns a list of the substrings as you can see in the output.

# The re module

```
import re
```

- Replace one or many matches with a string:

re.sub(r"regex", new, string)

```
re.sub(r"yellow", "nice", "I have a yellow car and a yellow house in a yellow neighborhood")
```

```
'I have a nice car and a nice house in a nice neighborhood'
```

# Supported metacharacters

| Metacharacter | Meaning |
|:---:|:---:|
| \d | Digit |

```python
re.findall(r"User\d", "The winners are: User9, UserN, User8")
```

```
['User9', 'User8']
```

| Metacharacter | Meaning |
|:---:|:---:|
| \D | Non-digit |

```python
re.findall(r"User\D", "The winners are: User9, UserN, User8")
```

```
['UserN']
```

13. Supported metacharacters
Let's look at the supported metacharacters. In the example, we want to find all matches of the patterns containing User followed by a number. We use backlash d to represent the digit. We get the following matches. Next, we find matches of the pattern containing User followed by a non-digit. In that case, we use backslash capital D obtaining the following match.

# Supported metacharacters

| Metacharacter | Meaning |
|:---:|:---:|
| \w | Word |

```
re.findall(r"User\w", "The winners are: User9, UserN, User8")
```

```
['User9', 'UserN', 'User8']
```

| Metacharacter | Meaning |
|:---:|:---:|
| \W | Non-word |

```
re.findall(r"\W\d", "This skirt is on sale, only $5 today!")
```

```
['$5']
```

14. Supported metacharacters
If we want to find all matches of the pattern containing User followed by any digit or normal character, we can use backlash w. We get all following matches. In the next example, we need to find the price in a string. We use backslash capital W to match the dollar sign followed by a digit obtaining the following output.

# Supported metacharacters

| Metacharacter | Meaning |
|:---:|:---:|
| \s | Whitespace |

```python
re.findall(r"Data\sScience", "I enjoy learning Data Science")
```

```
['Data Science']
```

| Metacharacter | Meaning |
|:---:|:---:|
| \S | Non-Whitespace |

```python
re.sub(r"ice\Scream", "ice cream", "I really like ice-cream")
```

```
'I really like ice cream'
```

15. Supported metacharacters
Finally, we use backslash s to specify the pattern Data whitespace science getting the following match. In the second example, we use backslash capital S to detect the matches of ice, followed by any non-space character, followed by cream and replace them with the word ice cream.

# Let's practice!

## REGULAR EXPRESSIONS IN PYTHON

# Repetitions

## REGULAR EXPRESSIONS IN PYTHON

**Maria Eugenia Inzaugarat**
Data Science

# Repeated characters

Validate the following string:

**password1234**

# Repeated characters

Validate the following string:



password1234

# Repeated characters

Validate the following string:

password1234

# Repeated characters

Validate the following string:



password1234

```python
import re
password = "password1234"
```

```python
re.search(r"\w\w\w\w\w\w\w\w\d\d\d\d", password)
```

```
<_sre.SRE_Match object; span=(0, 12), match='password1234'>
```

6. Repeated characters
**But this seems cumbersome for longer regex.** Instead, we can use quantifiers to save this situation. A quantifier is a metacharacter that specifies how many times a character located to its left needs to be matched. In our example, we specify that backslash w is repeated 8 times and backslash d four times. And we get a match as seen in the output.

# Repeated characters

Validate the following string:

password1234

```python
import re
password = "password1234"
```

```python
re.search(r"\w{8}\d{4}", password)
```

```
<_sre.SRE_Match object; span=(0, 12), match='password1234'>
```

**Quantifiers:**

A metacharacter that tells the regex engine how many times to match a character immediately to its left.

# Quantifiers

- Once or more: +

```
text = "Date of start: 4-3. Date of registration: 10-04."
```

```
re.findall(r"        ", text)
```

# Quantifiers

- Once or more: +

```python
text = "Date of start: 4-3. Date of registration: 10-04."
```

```python
re.findall(r"\d+-    ", text)
```

# Quantifiers

- Once or more: +

```
text = "Date of start: 4-3. Date of registration: 10-04."
```

```
re.findall(r"\d+-\d+", text)
```

```
['4-3', '10-04']
```

9. Quantifiers
And again a digit, backslash d. It should appear once or more times, so again we use the plus quantifier. We get the two matches that we expected as seen in the output.

# Quantifiers

- Zero times or more: *

```
my_string = "The concert was amazing! @ameli!a @joh&&n @mary90"

re.findall(r"@\w+\W*\w+", my_string)
```

```
['@ameli!a', '@joh&&n', '@mary90']
```

10. Quantifiers
To indicate that a character should appear zero or more times, we can use the star metacharacter. In this example, we have the following string. We want to find all mentions of users, which start with an @. We notice that they could contain or not contain a non-word character in the middle. So, we construct our regex as seen in the code: an @, followed by backslash w plus, backslash capital w with start to indicate a non-word character zero or more times, then backslash w plus. And we get the matches seen in the output.

# Quantifiers

- Zero times or once: ?

```
text = "The color of this image is amazing. However, the colour blue could be brighter."
re.findall(r"colou?r", text)
```

```
['color', 'colour']
```

11. Quantifiers

Another helpful quantifier is the question mark. It indicates that a character should appear zero times or once. In the example, we need to find matches for the word color, which has spelling variations. So, our regex is c, o, l, o, u which can appear once or zero times, so we add the question mark after it, then r. With this regex, we get the two matches wanted.

# Quantifiers

- <mark>n times at least,</mark> m times at most : `{n, m}`

```python
phone_number = "John: 1-966-847-3131 Michelle: 54-908-42-42424"
```

```python
re.findall(r"                                        ", phone_number)
```

# Quantifiers

- n times at least, m times at most : `{n, m}`

```
phone_number = "John: 1-966-847-3131 Michelle: 54-908-42-42424"
```

```
re.findall(r"\d{1,2}-                    ", phone_number)
```

13. Quantifiers
As we can see, we can have a digit, once or twice, then slash.

14. Quantifiers
Then, again a digit, three times. Then, a slash.

# Quantifiers

- n times at least, m times at most : `{n,  m}`

```
phone_number = "John: 1-966-847-3131 Michelle: 54-908-42-42424"
```

```
re.findall(r"\d{1,2}-\d{3}-              ", phone_number)
```

15. Quantifiers
Then, a digit twice or three times. Then, slash. And finally a digit again. We indicate that this last digit need to appear at least four times leaving the second argument in black. The regex engine returns two matches as expected.

# Quantifiers

- n times at least, m times at most : {n, m}

```python
phone_number = "John: 1-966-847-3131 Michelle: 54-908-42-42424"
```

```python
re.findall(r"\d{1,2}-\d{3}-\d{2,3}-\d{4,}", phone_number)
```

```
['1-966-847-3131', '54-908-42-42424']
```

16. Quantifiers
It's very important to remember that the quantifier applies only to the character immediately to the left. As an example, in the regex apple plus, the plus applies only to the letter e and not to the entire word.

# Quantifiers

- *Immediately to the left*
  - `r"apple+"` : **+** applies to e and not to apple

16. Quantifiers
It's very important to remember that the quantifier applies only to the character immediately to the left. As an example, in the regex apple plus, the plus applies only to the letter e and not to the entire word.

# Let's practice!

datacamp

# Regex metacharacters

## REGULAR EXPRESSIONS IN PYTHON

**Maria Eugenia Inzaugarat**

Data Scientist

# Looking for patterns

re.search(r"regex", string)

re.match(r"regex", string)

```
re.search(r"\d{4}", "4506 people attend the show")
```

```
<_sre.SRE_Match object; span=(0, 4), match='4506'>
```

```
re.match(r"\d{4}", "4506 people attend the show")
```

```
<_sre.SRE_Match object; span=(0, 4), match='4506'>
```

```
re.search(r"\d+", "Yesterday, I saw 3 shows")
```

```
<_sre.SRE_Match object; span=(17, 18), match='3'>
```

```
re.match(r"\d+","Yesterday, I saw 3 shows")
```

```
None
```

2. Looking for patterns
Let's first look at two methods of the re module: dot search and dot match. As you can see, they have the same syntax and are used to find a match.. Both methods return an object with the match found. The difference is that dot match is anchored at the **beginning of the string.** In this case, dot search finds a match, but dot match does not. This is because the first characters do not match the regex.

# Special characters

- <mark>Match any character</mark> (except newline): .

www.domain.com

```
my_links = "Just check out this link: www.amazingpics.com. It has amazing photos!"
re.findall(r"www  com", my_links)
```

3. Special characters
The dot metacharacter matches any character. In the example code, we need to match links in the string. We know a link starts with w, w, w and ends with c, o, m. So we first write this in our regex.

4. Special characters
We don't know how many character are in between. So we indicate that we want any character, a dot, once or more times, adding the plus. We can see in the output that we get our match.

# Special characters

- Match any character (except newline): .

www.**domain**.com

```python
my_links = "Just check out this link: www.amazingpics.com. It has amazing photos!"
re.findall(r"www.+com", my_links)
```

```
['www.amazingpics.com']
```

The dot . metacharacter matches any character.

# Special characters

- Start of the string: `^`

```
my_string = "the 80s music was much better that the 90s"
```

```
re.findall(r"the\s\d+s", my_string)
```

```
['the 80s', 'the 90s']
```

```
re.findall(r"^the\s\d+s", my_string)
```

5. Special characters
The circumflex anchors the regex to the start of a string. In the example, we find the pattern starting with t, followed by h, e, whitespace, two digits and ending with s. The method finds the following two matches. Now, we add the anchor metacharacter. **We receive only one match. The one that appears at the beginning of the string.**

```
['the 80s']
```

# Special characters

- End of the string: $

```
my_string = "the 80s music hits were much better that the 90s"
```

```
re.findall(r"the\s\d+s$", my_string)
```

```
['the 90s']
```

6. Special characters
On the contrary, the dollar sign anchors the regex to the **end of the string.**
If we use it in the previous example, we get the match that appears at the
end of the string.

# Special characters

- Escape special characters: \

```python
my_string = "I love the music of Mr.Go. However, the sound was too loud."
```

```python
print(re.split(r".\s", my_string))
```

```
['', 'lov', 'th', 'musi', 'o', 'Mr.Go', 'However', 'th', 'soun', 'wa', 'to', 'loud.']
```

```python
print(re.split(r"\.\s", my_string))
```

```
['I love the music of Mr.Go', 'However, the sound was too loud.']
```

# OR operator

- Character: |

```
my_string = "Elephants are the world's largest land animal! I would love to see an elephant one day"
```

```
re.findall(r"Elephant|elephant", my_string)
```

```
['Elephant', 'elephant']
```

8. OR operator
In the example code, we want to match the word elephant. However, we see that it's written with capital E or lower e. In that case, we use the vertical bar. In this way, we indicate that we want to match one variant OR the other obtaining both elephant-matches.

# OR operator

- Set of characters: `[ ]`

```python
my_string = "Yesterday I spent my afternoon with my friends: MaryJohn2 Clary3"
```

```python
re.findall(r"[a-zA-Z]+\d", my_string)
```

```
['MaryJohn2', 'Clary3']
```

9. OR operator
Square brackets also represent the OR operand. Inside them, we can specify optional characters to match. Look at the example. We want to find a pattern that contains lowercase or uppercase letter followed by a digit. To do so, we can use the square brackets. Inside them, we will use lowercase a dash lowercase z to specify any lowercase letter. Then, uppercase a dash uppercase z to indicate any uppercase letter. Then the plus. Then backslash d. Thus, we get the following matches.

10. OR operator
In the following string, we want to replace the non-word characters by whitespace. We specify optional characters inside the square brackets. The engine searches for one or the other. When it finds a match, it replaces the character by a whitespace, getting the following output.

# OR operator

- Set of characters: `[ ]`

```python
my_string = "My&name&is#John Smith. I%live$in#London."
```

```python
re.sub(r"[#$%&]", " ", my_string)
```

```
'My name is John Smith. I live in London.'
```

# OR operand

- Set of characters: `[  ]`
  - `^` transforms the expression to negative

```python
my_links = "Bad website: www.99.com. Favorite site: www.hola.com"
re.findall(r"www[^0-9]+com", my_links)
```

```
['www.hola.com']
```

# Let's practice!

REGULAR EXPRESSIONS IN PYTHON

# Greedy vs. non-greedy matching

## REGULAR EXPRESSIONS IN PYTHON

**Maria Eugenia Inzaugarat**
Data Scientist

# Greedy vs. non-greedy matching

- Two types of matching methods:
  - Greedy

  - Non-greedy or lazy

- Standard quantifiers are greedy by default: `*`, `+`, `?`, `{num, num}`

2. Greedy vs. non-greedy matching
There are two types of matching methods: greedy and non-greedy (also called lazy) operators.
The quantifiers that you have been learning until now (which are called standard quantifiers)
are greedy by default.

# Greedy matching

- **Greedy:** match as many characters as possible

- Return the *longest match*

```python
import re
re.match(r"\d+", "12345bcada")
```

```
<_sre.SRE_Match object; span=(0, 5), match='12345'>
```

\d+          \d+

12345abcde   12345abcde   12345abcde

# Greedy matching

- Backtracks when too many character matched

- Gives up characters one at a time

```
import re
re.match(r".*hello", "xhelloxxxxx")
```

```
<_sre.SRE_Match object; span=(0, 6), match='xhello'>
```

.*  →  xhelloxxxxx

.*h  →  xhelloxxxxx

.*h  →  xhelloxxxxx  »»»»  xhelloxxxxx  »»»»  xhelloxxxxx

.*hello

# Non-greedy matching

- **Lazy:** match as few characters as needed

- Returns *the shortest match*

- Append ? to greedy quantifiers

```
import re
re.match(r"\d+?", "12345bcada")
```

```
<_sre.SRE_Match object; span=(0, 1), match='1'>
```

\d+

12345abcde

\d+?

12345abcde

# Non-greedy matching

- Backtracks when too few characters matched

- Expands characters one a time

```python
import re
re.match(r".*?hello", "xhelloxxxxx")
```

```
<_sre.SRE_Match object; span=(0, 6), match='xhello'>
```

# Let's practice!

REGULAR EXPRESSIONS IN PYTHON