# Regular expressions



#### Class outline:

- Declarative languages
- Regular expression syntax
- Regular expressions in Python

# Declarative languages

# Declarative programming

#### In imperative languages:

- A "program" is a description of computational processes
- The interpreter carries out execution/evaluation rules

#### In declarative languages:

- A "program" is a description of the desired result
- The interpreter figures out how to generate the result
- Examples:
  - Regular expressions: Good (?:morning|evening)
  - Backus-Naur Form:

```
?calc_expr: NUMBER | calc_op
calc_op: "(" OPERATOR calc_expr* ")"
OPERATOR: "+" | "-" | "*" | "/"
```

# Domain-specific languages

Many declarative languages are **domain-specific**: they are designed to tackle problems in a particular domain, instead of being general purpose multi-domain programming languages.

Language	Domain
Regular expressions	Pattern-matching strings
Backus-Naur Form	Parsing strings into parse trees
SQL	Querying and modifying database tables
HTML	Describing the semantic structure of webpage content
CSS	Styling webpages based on selectors
Prolog	Describes and queries logical relations

# Regular expressions

# Pattern matching

Pattern matching in strings is a common problem in computer programming.

An imperative approach:

```
def is_email_address(str):
    parts = str.split('@')
    if len(parts) != 2:
        return False
    domain_parts = parts[1].split('.')
    return len(domain_parts) >= 2 and len(domain_parts[-1]) == 3
```

# Pattern matching

Pattern matching in strings is a common problem in computer programming.

An imperative approach:

```
def is_email_address(str):
    parts = str.split('@')
    if len(parts) != 2:
        return False
    domain_parts = parts[1].split('.')
    return len(domain_parts) >= 2 and len(domain_parts[-1]) == 3
```

An equivalent regular expression:

```
(.+)@(.+)\.(.{3})
```

With regular expressions, a programmer can just describe the pattern using a common syntax, and a regular expression engine figures out how to do the pattern matching for them.

# Matching exact strings

To match an exact string that has no special characters, just use the string:

Berkeley, CA 94720



Fully matched by: Berkeley, CA 94720

But if the matched string contains special characters, they must be escaped using a backslash.

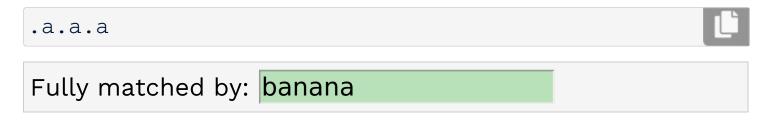
\(1\+3\)



Fully matched by: (1+3)

#### The dot

The . character matches any single character that is not a new line.



It's typically better to match a more specific range of characters, however...

# **Character classes**

<b>Pattern</b>	Description	Example	Fully Matched By
	Denotes a character class. Matches characters in a set (including ranges of characters like 0-9).  Use [^] to match characters outside a set.	<pre>[top] [h-p]</pre>	j
	Matches any character other than the newline character.	1.	1?
\d	Matches any digit character. Equivalent to [0-9]. \D matches the inverse (all non-digit characters).	\d\d	12
\W	Matches any word character. Equivalent to [A-Za-z0-9_]. \W matches the inverse.	\d\w	4Z
\s	Matches any whitespace character: spaces, tabs, or line breaks. \S matches the	\d\s\w	9 a

# Quantifiers

These indicate how many of a character/character class to match.

Pattern	<b>Description</b>	Example	Fully Matched By
*	Matches 0 or more of the previous pattern.	a*	aaa
+	Matches 1 or more of the previous pattern.	lo+l	lool
?	Matches 0 or 1 of the previous pattern.	lo?l	lol
{}	Used like {Min, Max}. Matches a quantity between	a{2}	aa
	Min and Max of the previous pattern.	a{2,}	aaaaaa
		a{2,4}	aaa

#### Combining patterns

Patterns P<sub>1</sub> and P<sub>2</sub> can be combined in various ways.

Combination	Description	Example	<b>Fully Matched By</b>
$P_1P_2$	A match for P <sub>1</sub> followed immediately by one for P <sub>2</sub> .	ab[.,]	ab,
P <sub>1</sub>   P <sub>2</sub>	Matches anything that either P1 or P2 does.	\d+ Inf	Inf
(P <sub>1</sub> )	Matches whatever P1 does. Parentheses group, just as in arithmetic expressions.	(<3)+	<3<3<3

#### **Anchors**

These don't match an actual character, they indicate the position where the surrounding pattern should be found.

Pattern	Description	Example	What parts match?
^	Matches the beginning of a string.	^aw+	aww aww
\$	Matches the end of a string.	\w+y\$	stay stay
\b	Matches a word boundary, the beginning or end of a word.	\w+e\b	broken bridge team

# Regular expressions in Python

# Support for regular expressions

Regular expressions are supported natively in many languages and tools.

Languages: Perl, ECMAScript, Java, Python, ...

Tools: Excel/Google Spreadsheets, SQL, BigQuery, VSCode, grep, ...

#### Raw strings

In normal Python strings, a backslash indicates an escape sequence, like \n for new line or \b for bell.

```
>>> print("I have\na newline in me.")
I have
a newline in me
```

But backslash has a special meaning in regular expressions. To make it easy to write regular expressions in Python strings, use raw strings by prefixing the string with an r:

```
pattern = r"\b[ab] +\b"
```

#### The re module

The re module provides many helpful functions.

Function	Description
<pre>re.search(pattern, string)</pre>	returns a Match object representing the first occurrence of pattern within string
<pre>re.fullmatch(pattern, string)</pre>	returns a Match object, requiring that pattern matches the entirety of string
<pre>re.match(pattern, string)</pre>	returns a Match object, requiring that string starts with a substring that matches pattern
<pre>re.findall(pattern, string)</pre>	returns a list of strings representing all matches of pattern within string, from left to right
<pre>re.sub(pattern, repl, string)</pre>	substitutes all matches of pattern within string with repl

# Match objects

The functions re.search, re.match, and re.fullmatch all take a string containing a regular expression and a string of text. They return either a Match object or, if there is no match, None.

re.search requires that the pattern exists somewhere in the string:

```
import re

re.search(r'-?\d+', '123 peeps')  # <re.Match object>
re.search(r'-?\d+', 'So many peeps')  # None
```

# Match objects

The functions re.search, re.match, and re.fullmatch all take a string containing a regular expression and a string of text. They return either a Match object or, if there is no match, None.

re.search requires that the pattern exists somewhere in the string:

```
import re

re.search(r'-?\d+', '123 peeps')  # <re.Match object>
re.search(r'-?\d+', 'So many peeps')  # None
```

Match objects are treated as true values, so you can use the result as a boolean:

```
bool(re.search(r'-?\d+', '123'))  # True
bool(re.search(r'-?\d+', 'So many peeps')) # False
```

# Inspecting a match

re.search returns a Match object representing the first occurrence of pattern within string.

```
title = "I Know Why the Caged Bird Sings"
re.search(r'Bird', title) #
```

Match objects carry information about what has been matched. The Match.group() method allows you to retrieve it.

```
x = "This string contains 35 characters."
mat = re.search(r'\d+', x)
mat.group(0) # 35
```

If there are parentheses in a patterns, each of the parenthesized groups will become groups in the match object.

```
x = "There were 12 pence in a shilling and 20 shillings in a pound."
mat = re.search(r'(\d+)[a-z\s]+(\d+)', x)

mat.group(0)
mat.group(1)
mat.group(2)
mat.groups()
```

If there are parentheses in a patterns, each of the parenthesized groups will become groups in the match object.

```
x = "There were 12 pence in a shilling and 20 shillings in a pound."
mat = re.search(r'(\d+)[a-z\s]+(\d+)', x)

mat.group(0)  # '12 pence in a shilling and 20'
mat.group(1)
mat.group(2)
mat.groups()
```

If there are parentheses in a patterns, each of the parenthesized groups will become groups in the match object.

```
x = "There were 12 pence in a shilling and 20 shillings in a pound."
mat = re.search(r'(\d+)[a-z\s]+(\d+)', x)

mat.group(0)  # '12 pence in a shilling and 20'
mat.group(1)  # 12
mat.group(2)
mat.groups()
```

If there are parentheses in a patterns, each of the parenthesized groups will become groups in the match object.

```
x = "There were 12 pence in a shilling and 20 shillings in a pound."
mat = re.search(r'(\d+)[a-z\s]+(\d+)', x)

mat.group(0)  # '12 pence in a shilling and 20'
mat.group(1)  # 12
mat.group(2)  # 20
mat.groups()
```

If there are parentheses in a patterns, each of the parenthesized groups will become groups in the match object.

```
x = "There were 12 pence in a shilling and 20 shillings in a pound."
mat = re.search(r'(\d+)[a-z\s]+(\d+)', x)

mat.group(0)  # '12 pence in a shilling and 20'
mat.group(1)  # 12
mat.group(2)  # 20
mat.groups()  # (12, 20)
```

# **Exercises**

[A-Za-z] {3}

Fully matched by: ?

- What's a valid input?
- What's an invalid input?

[A-Za-z] {3}

Fully matched by: ?

- What's a valid input? AUS, aus
- What's an invalid input? australia, au

\d{4} - \d{2} - \d{2}

Fully matched by: ?

- What's a valid input?
- What's an invalid input?

\d{4} - \d{2} - \d{2}

Fully matched by: ?

- What's a valid input? 2020-03-13
- What's an invalid input? 2020/03/13, 03-13-2020

$$[a-z0-9._%+-]+@[a-z0-9.-]+\.[a-z]{2,}$$
\$



Fully matched by: ?

- What's a valid input?
- What's an invalid input?

$$[a-z0-9._%+-]+@[a-z0-9.-]+\.[a-z]{2,}$$
\$



Fully matched by: ?

- What's a valid input? someone@someplace.org
- What's an invalid input? someone@mod%cloth.co

#### **Exercise: Stocks**

Make a regular expression to match any tweet talking about GME stock.

```
import re

def match_gme(tweet):
    """
    >>> match_gme('GME')
    True
    >>> match_gme('yooo buy GME right now!')
    True
    >>> match_gme('#HUGME')
    False
    >>> match_gme('#HUGMEHARDER')
    False
    """
    return bool(re.search(_____, tweet))
```

# Tips

- When learning, use sites like regexr.com
- Get used to referencing a regular expressions cheat sheet

#### $\triangle$ A word of caution $\triangle$

Regular expressions can be very useful. However:

- Very long regular expressions can be difficult for other programmers to read and modify.
   See also: Write-only
- Since regular expressions are declarative, it's not always clear how efficiently they'll be processed. Some processing can be so time-consuming, it can take down a server.
- Regular expressions can't parse everything! Don't write an HTML parser with regular expressions.