

PATTERN MATCHING WITH REGULAR EXPRESSIONS

CS 3080: Python Programming



University of Colorado
Colorado Springs

Regular expressions

- You may be familiar with searching for text by pressing ctrl-F and typing in the words you're looking for.
- Regular expressions go one step further: They **allow you to specify a pattern of text to search for.**

“Knowing [regular expressions] can mean the difference between solving a problem in 3 steps and solving it in 3,000 steps. When you’re a nerd, you forget that the problems you solve with a couple keystrokes can take other people days of tedious, error-prone work to slog through.”

- Cory Doctorow

Time to code: Finding Patterns of Text Without Regular Expressions

- You want to find a phone number in a string. You know the pattern: three numbers, a hyphen, three numbers, a hyphen, and four numbers.
- Here's an example: 415-555-4242.
- Let's write a function `isPhoneNumber()` to check whether a string matches this pattern, returning either `True` or `False`

Finding Patterns of Text with Regular Expressions

- What if the number is 415.555.4242, (415) 555-4242, 415-555-4242 x99? The `isPhoneNumber()` function would fail to validate them. You could add yet more code for these additional patterns, but there is an easier way.

Finding Patterns of Text with Regular Expressions

- What if the number is 415.555.4242, (415) 555-4242, 415-555-4242 x99? The `isPhoneNumber()` function would fail to validate them. You could add yet more code for these additional patterns, but there is an easier way.
- **Regular expressions**, called *regexes* for short, are descriptions for a pattern of text.
- For example, a `\d` in a regex stands for a digit character (any single numeral 0 to 9). So, the `isPhoneNumber()` function as a regex could be as:
 - `\d\d\d-\d\d\d-\d\d\d\d`
- Or even more sophisticated:
 - `\d{3}-\d{3}-\d{4}`

Finding Patterns of Text with Regular Expressions

- What if the number is 415.555.4242, (415) 555-4242, 415-555-4242 x99? The `isPhoneNumber()` function would fail to validate them. You could add yet more code for these additional patterns, but there is an easier way.
- **Regular expressions**, called *regexes* for short, are descriptions for a pattern of text.
- For example, a `\d` in a regex stands for a digit character (any single numeral 0 to 9). So, the `isPhoneNumber()` function as a regex could be as:
 - `\d\d\d-\d\d\d-\d\d\d\d`
- Or even more sophisticated:
 - `\d{3}-\d{3}-\d{4}`

Complexity of a regex is $O(n)$
where n is the length of the text to search.

Creating a regex object

(1) Pass desired pattern to `re.compile()`, and store the resulting `Regex` object in `phoneNumRegex`

(2) Then call `search()` on `phoneNumRegex`, and pass `search()` the string we want to search. The result of the search is stored in `mo` (a generic variable name to use for `Match` objects)

(3) Call `group()` on `mo` to return the match.
`mo.group()` displays the whole match

Review of Regular Expression Matching

- While there are several steps to using regular expressions in Python, each step is fairly simple:
 1. *Import the regex module with `import re`.*
 2. *Create a Regex object with the `re.compile()` function. (Remember to use a raw string.)*
 3. *Pass the string you want to search into the Regex object's `search()` method. This returns a Match object or None if nothing is found.*
 4. *Call the Match object's `group()` method to return a string of the actual matched text.*

Grouping with Parenthesis

- Adding **parentheses** will create groups in the regex: `(\d\d\d)-(\d\d\d-\d\d\d\d)`.
- Parentheses have a special meaning in regular expressions, but what do you do if you need to match a parenthesis in your text?
 - *For example, area code in parenthesis*

```
phoneNumRegex = re.compile(r'(\d\d\d) (\d\d\d-\d\d\d\d)')
```

- The `\(` and `\)` escape characters in the raw string passed to `re.compile()` will match actual parenthesis characters.

Matching Multiple Patterns with Pipe

- The | character is called a pipe. You can use it anywhere you want to match one of many expressions.
 - *For example, the regular expression `r'Batman|Tina Fey'` will match either 'Batman' or 'Tina Fey'.*
 - *When both Batman and Tina Fey occur in the searched string, **the first occurrence of matching text** will be returned as the Match object.*
- You can also use the pipe to match one of several patterns as part of your regex.
- *If you need to match an actual pipe character, escape it with a backslash, like `\|`. The same happens with all the following characters in the next slides.*

Optional Matching with Question Mark

- The ? character flags the group that precedes it as an optional part of the pattern.
- The (wo)? Part of the regular expression means that the pattern wo is an optional group
- The regex will match text that has zero instances or one instance of wo in it
- This is why the regex matches both Batwoman and Batman

Matching Zero or More with the Star

- The * (called the star or asterisk) means “match zero or more”—the group that precedes the star can occur any number of times in the text. It can be completely absent or repeated over and over again.
- For 'Batman' , the (wo)* part of the regex matches zero instances of wo in the string;
- For 'Batwoman' , the (wo)* matches one instance of wo ;
- For 'Batwowowowoman' , (wo)* matches four instances of wo .

Matching Specific Repetitions with Curly Brackets

- If you have a group that you want to repeat a specific number of times, follow the group in your regex with a number in curly brackets.
 - *For example, the regex `(Ha){3}` will match the string 'HaHaHa', but it will not match 'HaHa'*
- These two regular expressions match identical patterns:
 - `(Ha){3}`
 - `(Ha)(Ha)(Ha)`
- But `(Ha){3}` returns only one group

```
print(mo.group(0))      # HaHaHa
print(mo.group(1))      # Ha
print(mo.group(2))      # IndexError: no such group
```

Matching Specific Repetitions with Curly Brackets

- Instead of one number, you can specify a range by writing a minimum, a comma, and a maximum in between the curly brackets.
 - *For example, these two regular expressions match identical patterns:*
 - *(Ha){3,5}*
 - *((Ha)(Ha)(Ha))|((Ha)(Ha)(Ha)(Ha))|((Ha)(Ha)(Ha)(Ha)(Ha))*
- You can also leave out the first or second number in the curly brackets to leave the minimum or maximum unbounded. For example,
 - *(Ha){3,} will match three or more instances of the (Ha) group,*
 - *(Ha){,5} will match zero to five instances*

Greedy and Nongreedy Matching

- Python's regular expressions are **greedy** by default, which means that in ambiguous situations they will match **the longest string possible**.
- The **nongreedy** version of the curly brackets, which matches the shortest string possible, has the closing curly bracket followed by a question mark.

findall() method

- While `search()` will return a `Match` object of the first matched text in the searched string, the **`findall()`** method will return the strings of every match in the searched string.
- If there are groups in the regular expression, then `findall()` will return a list of tuples. Each tuple represents a found match, and its items are the matched strings for each group in the regex.

Character classes

- The character class `\d` is shorthand for the regular expression `(0|1|2|3|4|5|6|7|8|9)`.
- The character class `[0-5]` will match only the numbers 0 to 5; this is much shorter than typing `(0|1|2|3|4|5)`.

Table 7-1: Shorthand Codes for Common Character Classes

Shorthand character class	Represents
<code>\d</code>	Any numeric digit from 0 to 9.
<code>\D</code>	Any character that is <i>not</i> a numeric digit from 0 to 9.
<code>\w</code>	Any letter, numeric digit, or the underscore character. (Think of this as matching “word” characters.)
<code>\W</code>	Any character that is <i>not</i> a letter, numeric digit, or the underscore character.
<code>\s</code>	Any space, tab, or newline character. (Think of this as matching “space” characters.)
<code>\S</code>	Any character that is <i>not</i> a space, tab, or newline.

Character class example

- The regular expression `\d+\s\w+` will match text that has one or more numeric digits (`\d+`), followed by a whitespace character (`\s`), followed by one or more letter/digit/underscore characters (`\w+`).
- The `findall()` method returns all matching strings of the regex pattern in a list.

Making Your Own Character Classes

- There are times when you want to match a set of characters but the shorthand character classes (`\d`, `\w`, `\s`, and so on) are too broad. You can create your own character classes using square brackets:
 - *[aeiouAEIOU], will match any vowel, both lowercase and uppercase*
 - *[a-zA-Z0-9], will match all lowercase letters, uppercase letters, and numbers*
 - *[a-z], will match all lowercase letters*
 - *etc*

Making Your Own Character Classes

- Note that inside the square brackets, the normal regular expression symbols are not interpreted as such.
- This means you **do not need to escape** the ., *, ?, or () characters with a preceding backslash. For example,
 - *[0-5.] will match digits 0 to 5 and a period.*
 - *You do not need to write it as [0-5\.]*
 - *If you want to match the backslash too, you have to scape it: [0-5\\.]*

Making Your Own Character Classes

- By placing a **caret character (^)** just after the character class's opening bracket, you can make a **negative character class**. A negative character class will match all the characters that are not in the character class.

The Caret and Dollar Sign Characters

- You can also use the **caret symbol (^)** at the start of a **regex** to indicate that a match must occur at the beginning of the searched text.
- Likewise, you can put a **dollar sign (\$)** at the end of the **regex** to indicate the string must end with this regex pattern.
- And you can use the **^ and \$ together** to indicate that the entire string must match the regex.

The Wildcard Character

- The . (or dot) character in a regular expression is called a **wildcard** and will match **any character** except for a newline.
- And to include the newline character too:

```
newlineRegex = re.compile(r'.at', re.DOTALL)
```

Matching Everything with Dot-Star

- The dot-star uses **greedy** mode: It will always try to match as much text as possible.
- To match any and all text in a **nongreedy** fashion, use the dot, star, and question mark (`.*`).
- Like with curly brackets in `(Ha){3,5}?`, the question mark tells Python to match in a nongreedy way.

Review of Regex Symbols 1

- The `?` matches zero or one of the preceding group.
- The `*` matches zero or more of the preceding group.
- The `+` matches one or more of the preceding group.
- The `{n}` matches exactly `n` of the preceding group.
- The `{n,}` matches `n` or more of the preceding group.
- The `{,m}` matches 0 to `m` of the preceding group.
- The `{n,m}` matches at least `n` and at most `m` of the preceding group.
- `{n,m}?` or `*?` or `+` performs a nongreedy match of the preceding group.

Review of Regex Symbols 2

- `^spam` means the string must begin with spam.
- `spam$` means the string must end with spam.
- The `.` matches any character, except newline characters.
- `\d`, `\w`, and `\s` match a digit, letter, or space character, respectively.
- `\D`, `\W`, and `\S` match anything except a digit, letter, or space character, respectively.
- `[abc]` matches any character between the brackets (such as a, b, or c).
- `[^abc]` matches any character that isn't between the brackets.

Case-Insensitive Matching

When you only care about matching the letters without worrying whether they are upper case or lower case.

Substituting Strings with the `sub()` Method

Regular expressions can not only find text patterns but can substitute new text in place of those patterns.

Substituting Strings with the `sub()` Method

- Sometimes you may need to use the matched text itself as part of the substitution.
- In the first argument to `sub()` , you can use `\1` , `\2` , `\3` , and so on, to mean “Enter the text of group 1 , 2 , 3 , and so on, in the substitution.”
- For example, say you want to censor the names of the secret agents by showing just the first letters of their names. To do this, you could use the regex `Agent (\w)\w*` and pass `r'\1****'` as the first argument to `sub()` . The `\1` in that string will be replaced by whatever text was matched by group 1 — that is, the `(\w)` group of the regular expression.

Managing Complex Regexes

```
phoneRegex = re.compile(r'((\d{3}|\d{3}\d{3})?(\s|-|\.)?\d{3}(\s|-|\.)\d{4}(\s*(ext|x|ext.)\s*\d{2,5})?)')
```

- Use the re.VERBOSE and triple-quote syntax

```
phoneRegex = re.compile(r'''(
    (\d{3}|\d{3}\d{3})?      # area code
    (\s|-|\.)?              # separator
    \d{3}                   # first 3 digits
    (\s|-|\.)               # separator
    \d{4}                   # last 4 digits
    (\s*(ext|x|ext.)\s*\d{2,5})? # extension
)''', re.VERBOSE)
```

Combining `re.IGNORECASE`, `re.DOTALL` , and `re.VERBOSE`

```
someRegexValue = re.compile('foo', re.IGNORECASE | re.DOTALL | re.VERBOSE)
```