String Processing and Pattern Matching

Informatics 1 for Biomedical Engineers **Tutor Session 4**

KTI, Knowledge Technologies Institute

9. November 2016

http://kti.tugraz.at/







Today's Topics

- Strings in python
- Pattern Matching
- Regular Expressions



Student Goals

- Learn about the characteristics of strings in python
- Learn how to process strings: find, join, replace,...
- Understand the basics of regular expressions and how to use them



Data types revisited

- Logic: boolean
- Numeric types: int, float
- Sequences: list, tuple
- Text Sequence: str
- . . .



Strings in python

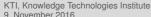
String variable

```
1 a = "Good_morning"
2 b = 'Good_evening'
```

Multi-line string variable

```
1 a = '''Good morning,
2 How are you feeling today?'''
```





Special characters

- How to define a string with line breaks/characters that break the syntax?
- → Escape characters

```
1 print('Hi,\n\'how\'uareuyou?')
```



Output:

```
Hi, 2 'how' are you?
```



Important Escape Characters¹

Special Character	Python Encoding
line break	\ n
tab	\ t
	\\
,	\', \''

Easy solution for string containing ' or ":

```
string = 'Hi, "how" are you?'
```



¹https://docs.python.org/3/reference/lexical_analysis.

html#string-and-bytes-literals



Working with Strings

- Concatenate, Format
- Substrings
- Contains
- Find
- Replace, Strip
- Split
- Join

Strings are immutable, always remember to store results in a (new) variable.



Working with Strings – Concat and Format

- Concatenate from separate parts
- Add variable values

```
# Concatenation
a = 'Part_1' + 'part_2.'

# Format with variable
field_of_study = 'Field_of_study:__{\}'.format('Biomedical_Engineering)
bpm = 'Blood_pressure:_systolic_{\}0\_mmHg,_diastolic_{\}1\_mmHg'
bpm_formatted = bpm.format(118, 70)

# convert to lower case
lower_case = a.lower()
```





Working with Strings – Substrings

- strings are sequences of characters
- define a start and end offset using square brackets: [0:2]
- important: start at index 0, end at len(string) 1
- leave out indices for entire start/end: [:]
- you can also count from the back: [:-2]

```
1 a = 'Hello_World.'
2 substring1 = a[2:4] # 11
3 substring2 = a[:] # Hello World.
4 substring3 = a[2:-2] # 110 Worl
```





Working with Strings – Contains

- check if one string contains another
- returns boolean value True or False

```
'Hel' in 'HellouWorld' # True
'hel' in 'HellouWorld' # False
```







Working with Strings – Find

- get the index at which a string occurs in another one
- returns -1 if the searched string was not found
- only returns the first occurrence!







Working with Strings - Replace

replace all occurrences of one string in another one

```
a = 'Hello_World.'

replace1 = a.replace('_', ';') # Hello; World.

replace2 = a.replace('1', ';') # He;; o Wor; d.

replace3 = a.replace('llo', 'x') # Hex World.
```





Working with Strings – Strip

remove whitespace characters at the beginning/end of a string

```
a = '_Hello_World._'

strip = a.strip() # 'Hello World.'
```







Working with Strings – Split

- break a string at the occurrence of a specified string
- returns a sequence of the parts, without the split characters
- always returns a string, even when the input is a number







Working with Strings – Join

- combine a sequence to a single string
- using a specified separator string

```
sequence = ['a', 'b', 'c']
joined1 = 'u'.join(sequence) # a b c

# join a sequence of characters, i.e. a string
separator = ','
joined2 = separator.join('Hello_World.') # H,e,l,l,o, ,W,o,r,l,d,.
```



Regular Expressions – Introduction

- What are regular expressions and what are they used for?
 - search occurrences of a given pattern
 - for search/replace implementation
- How do they work?
 - define the pattern you are looking for
 - e.g. "The letter A followed by two lower-case characters and one digit"
 - match a string/text based on this pattern



Regular Expressions – Sample Expressions²

- digits: [0-9], \d
- whitespace: \s
- wildcard (arbitrary character): .
- Encoding for "The letter A followed by two lower-case characters and one digit":
 A [a = 3 [a = 3] d

$$A[a-z][a-z] d$$

²https://docs.python.org/3/library/re.html



Regular Expressions – Email Addresses

The patterns can be arbitrarily complex. Example: Email validation according to RFC standard

```
(?:[a-z0-9!#$%&'*+/=?^_'{|}~-]+(?:\.[a-z0-9!#$%&'*+/=?^_'{|}~-]+)*
         x0bx0cx0e-x7f]*")@(?:(?:[a-z0-9](?:[a-z0-9-]*[a-z0-9])?.)
         +[a-z0-9](?:[a-z0-9-]*[a-z0-9])?[\[(?:(?:25[0-5]]2[0-4][0-9]][01]?
5
         [0-9][0-9]?) \setminus .) \{3\} (?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?|[a-z0-9-]*
6
         [a-z0-9]:(?:[\x01-\x08\x0b\x0c\x0e-\x1f\x21-\x5a\x53-\x7f]]
         [x01-x09x0bx0cx0e-x7f]+)
```



Regular Expressions – Groups

- Which pattern was matched by which part of the text? How to retrieve and address that part?
- \blacksquare \rightarrow use groups which are filled with the matching result
- defined with (<expression>)
- E.g. ([A-Z] [a-z])11 for "Hello World": group 0 is "He"
- also used to describe frequency of partial patterns, e.g. (ab)+|(cde)



Regular Expressions in python

- using the regular expressions module³
- general usage:
 - re.search(<pattern>, <string>)
 returns a match object which contains the whole match
 and groups
 - re.findall(<pattern>, <string>)
 returns a list of the matched strings

³https://docs.python.org/3/library/re.html



Regular Expressions in python

```
# import the regular expressions module
             import re
 3
            # re.search(pattern, string)
5
6
7
8
            match1 = re.search('1', 'Hello World.')
             # the entire match is saved as group 0
            print(match1.group(0)) # 1
9
            # all matches as a list
10
            match2 = re.findall('l[a-z]', 'Hello, World.')
11
            print(match2) # ['11', 'ld']
```





Regular Expressions in python – Groups

- address groups in the match object using the group index
- the 0th group is the entire matched string

```
# find parts of a phone number
match3 = re.search('(\+\d+)-(\d+)', '+43-680-1234567')
print('entire_number_' + match3.group(0)) # entire number +43-680-12
print('country_code_' + match3.group(1)) # country code +43
print('operator_code_' + match3.group(2)) # operator code 680
print('phone_number_' + match3.group(3)) # phone number 1234567
```



Complex example - Hangman Part 1

- 1. Search for a word list file⁴
- 2. Optionally clean that list or shorten it
- 3. Save that word list as a file and load it into python
- 4. Store all unique words. Hint: use a set⁵

⁴e.g. https://github.com/dwyl/english-words

⁵https://docs.python.org/3/tutorial/datastructures.html#sets



Complex example - Hangman Part 1

1. Randomly select one of the words. This will be the solution for a round of hangman

```
import random

# get a random sample from the set.

# Returns a list, so we choose the first (and only) element

selected_word = random.sample(<set>, 1)[0]
```



Complex example - Hangman Part 3

- 1. Using console input, the player can guess *one* letter at a time
- 2. Check if the random word contains that letter
- 3. Track progress:
 - how many guesses does the user have left?7 guesses are allowed
 - which characters were guessed correctly?
- 4. Notify the user when she has won/lost



```
# Load words into a set
word_list = open('words.txt', 'r')

words = []

for line in word_list:
    words.append(line.strip())

word_list.close()
```





See instruction slide

```
# Select a random word

import random
selected_word = random.choice(words)

# make it lower case so we don't have to worry about cases
selected_word = selected_word.lower()
```





```
# Set up game tracking variables

correct_letters = [False] * len(selected_word)

guessed_letters = []

current_word = ['_'] * len(selected_word)

wrong_letters_count = 0
```





```
# Ask player for character input
     while wrong_letters_count < 7 and False in correct_letters:</pre>
         # ask the user to enter a letter
         letter = input('Guess_a_letter:_')
 5
         letter = letter.lower()
 6
         # only one letter at a time!
         while (len(letter) is not 1):
8
             letter = input('Guessioneiletterionly:..')
9
             letter = letter.lower()
10
         # insert [Step 5] here
11
12
      # the word has not been guessed but all guesses are used up - the player le
13
     if False in correct letters:
14
         print('Sorry...vou..lost.')
15
     # everything was guessed correctly
16
     else:
17
         print('Yay, you won!')
  KTI, Knowledge Technologies Institute
  9 November 2016
```



```
# has this letter already been guessed before? Only continue if it hasn't
        if letter not in guessed_letters:
            # the word contain the player's letter!
            if letter in selected word:
5
                # insert [Step 6] here
6
            else:
8
                # that means one attempt less available...
                wrong_letters_count += 1
10
                print('Sorry, word does not contain ' + letter + '. You have
11
                    + str(7 - wrong_letters_count + 1) + '__guesses__remaining.')
12
13
        # this letter cannot be used anymore
14
        guessed letters.append(letter)
```

```
# now where do we find the letter?
         # We're not using find() because it only finds the first occurrence
         for i in range(len(selected_word)):
            if selected word[i] == letter:
5
                # this letter was guessed correctly.
6
                # so let's set that position to True
                correct letters[i] = True
8
                # show the player the current status by showing all
10
                # correct letters in the word
11
                current word[i] = selected word[i]
12
        # the known letters are in a list, so to print them correctly,
13
         # they need to be joined to a string
14
         print('This_is_iwhat_you_know_so_far:_' + ',' .join(current_word))
```



Student task – find a DNA pattern in a DNA string

- Load the DNA of baker's yeast⁶ and the genome AXL2⁷ from a file. Write a function for this!
- normalise the data: remove all whitespace characters ''
- Find some shorter base sequences, e.g. "AGT", "GTCC",... Not only the first occurrence, but all of them!
- Does baker's yeast have the AXL2 gene?

⁶http://www.ncbi.nlm.nih.gov/genbank/samplerecord
7-

⁷http://www.yeastgenome.org/locus/S000001402/overview (both preprocessed in the folder for this unit)

5 6

8

10

11 12

13 14



```
# Define a function to read a DNA file and return it as a string
def get_dna_sequence(file_name):
   dna_sequence = ""
   dna_file = open(file_name, "r")
   for line in dna_file:
       dna_sequence += line
   dna_sequence = dna_sequence.replace("", "")
   dna_sequence = dna_sequence.replace("\n", "")
   dna file.close()
   return dna sequence
```





```
# Load the baker's yeast file
# check for some simple base sequences
yeast_dna = get_dna_sequence("bakers_yeast.txt")

find_agt = yeast_dna.find("AGT")
print(find_agt) # 72
```





```
# That only gives us the first occurrence
     # - let's make it recursive and look at substrings!
     matches = \Pi
     index = 0
     sequence_to_find = "AGT"
6
     while True:
         new_index = yeast_dna[index:].find(sequence_to_find)
8
         if new index is -1:
            break
10
        else:
11
            index += new index + 1
12
            matches.append(index)
     print(matches) # [73, 90, 110, ..., 4926]
13
14
     # analogously for other sequences
```



2



```
# Load the AXL2 genome and check if baker's yeast has it
genome_dna = get_dna_sequence("AXL2_genomic.txt")
print(yeast_dna.find(genome_dna)) # True
```

