

Python Tutorial 2

LING-381-Language Technology and LLMs

Instructor: Hakyung Sung

September 11, 2025

Review: Tutorial 1

How to learn to ride a bike is ...



Review: Tutorial 1

to actually ride a bike.



Review: Basic value types

Python has several basic types of values. The most common include:

- **Strings** (text)

Review: Basic value types

Python has several basic types of values. The most common include:

- **Strings** (text)
- **Integers** (whole numbers)

Review: Basic value types

Python has several basic types of values. The most common include:

- **Strings** (text)
- **Integers** (whole numbers)
- **Floats** (decimal numbers)

Review: Basic value types

Python has several basic types of values. The most common include:

- **Strings** (text)
- **Integers** (whole numbers)
- **Floats** (decimal numbers)
- **Booleans** (*True* and *False*)

Review: Basic value types

Python has several basic types of values. The most common include:

- **Strings** (text)
- **Integers** (whole numbers)
- **Floats** (decimal numbers)
- **Booleans** (*True* and *False*)
- **None** (represents “no value”)

Review: Store values

Values

- Store values of any type.

```
1 a = "thing is a string"  
2 b = 9  
3 c = 3.2
```

```
1 print (b+c)
```

```
↵ 12.2
```

```
1 print (a+"!")
```

```
↵ thing is a string!
```

Review: Functions

Functions

```
1 print("Hi!")
```

```
Hi!
```

```
1 len("abc")
```

```
3
```

```
1 str(16)
```

```
'16'
```

```
1 int(3.9)
```

```
3
```

```
1 float("2.5")
```

```
2.5
```

```
1 type("abc") # you can check the type of the variables using this function
```

```
str
```

```
1 type(16)
```

```
int
```

Methods

```
1 sample = "This is a STRING"
```

```
1 sample.lower()
```

```
↔ 'this is a string'
```

```
1 sample.lower().split(".")
```

```
↔ ['this', 'is', 'a', 'string']
```

```
1 sample = "This, is, a, STRING"
```

```
▶ 1 # How can we modify this line?
```

```
↔ ['this', 'is', 'a', 'string']
```

Review: Membership test

Membership test

- Use the `in` operator to check if a substring exists:

```
1 word = "awesome"
2
3 if "a" in word:
4 | print ("Contains 'a'!")
5 else:
6 | print("No 'a' found.")
```

↩ Contains 'a'!

Tip: use `elif` for checking additional conditions

- `elif` = "else if"
- Used to check additional conditions after an initial `if` statement.
- only checked if the first condition was `False`

```
1 word = "awesome"
2
3 # we now want to check 'e' as well.
4 if "a" in word:
5 | print ("Contains 'a'")
6 elif "e" in word:
7 | print ("Contains 'e'")
8 else:
9 | print("No 'a' or 'e' found.")
```

↩ Contains 'a'

Review: for and while loop

Loops

```
1 words = ["read", "code", "repeat"]
2
3 for w in words:
4     print(w)
```

read
code
repeat

```
1 for w in words:
2     if "e" in w:
3         print(w)
```

read
code
repeat

```
1 count = 0
2 while count < 3:
3     print(count)
4     count += 1
```

0
1
2

```
1 while w in words:
2     print(w)
```

Review: Tuple

3. tuples, dictionaries, functions, classes, save into files

```
1 # Create and concatenate
2 t = (1, 2)
3 t = t + (3, 4)
4 print(t)
```

⇒ (1, 2, 3, 4)

```
1 # Indexing and slicing
2 print(t[0])
3 print(t[-1])
4 print(t[1:3])
```

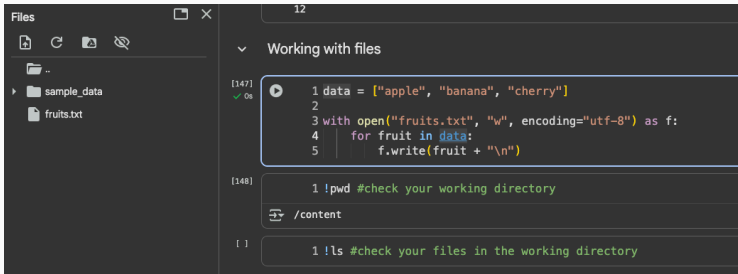
⇒ 1
4
(2, 3)

```
1 # Unpacking
2 x, y, *rest = t
3 print(x, y, rest)
4 # Explanation:
5 # x gets the first value (1), y gets the second value (2),
6 # and *rest collects all the remaining values into a list → [3, 4].
```

⇒ 1 2 [3, 4]

```
1 *begin, last = t
2 print(begin, last)
3
4 first, *middle, last = t
5 print(first, middle, last)
```

Review: Working with files



The screenshot shows a terminal window with a file explorer on the left and a code editor on the right. The file explorer shows a directory structure with 'sample_data' and 'fruits.txt'. The code editor shows a Python script that writes a list of fruits to a file, followed by shell commands to check the current directory and list files.

```
Files
├── ..
├── sample_data
└── fruits.txt
```

```
12
Working with files
[147] 1 data = ["apple", "banana", "cherry"]
      2
      3 with open("fruits.txt", "w", encoding="utf-8") as f:
      4     for fruit in data:
      5         f.write(fruit + "\n")

[148] 1 !pwd #check your working directory

      ↵ /content

[ ] 1 !ls #check your files in the working directory
```

- We often collect text data from the web.

- We often collect text data from the web.
- Raw text can be **messy**: strange symbols, broken characters, odd formatting.

Tutorial 2

- We often collect text data from the web.
- Raw text can be **messy**: strange symbols, broken characters, odd formatting.
- If we put messy text into a *DataFrame*, it may look confusing or inconsistent.

Tutorial 2

- We often collect text data from the web.
- Raw text can be **messy**: strange symbols, broken characters, odd formatting.
- If we put messy text into a *DataFrame*, it may look confusing or inconsistent.
- **First step: clean the text** so that we can work with it safely.

Tutorial 2

- We often collect text data from the web.
- Raw text can be **messy**: strange symbols, broken characters, odd formatting.
- If we put messy text into a *DataFrame*, it may look confusing or inconsistent.
- **First step: clean the text** so that we can work with it safely.
- After cleaning, we can move on to:
 - **Tokenization**

Tutorial 2

- We often collect text data from the web.
- Raw text can be **messy**: strange symbols, broken characters, odd formatting.
- If we put messy text into a *DataFrame*, it may look confusing or inconsistent.
- **First step: clean the text** so that we can work with it safely.
- After cleaning, we can move on to:
 - **Tokenization**
 - **Lemmatization**

- **Tokenization:** Converting raw text into tokens (words, punctuation, numbers). Essential for breaking text into analyzable units.

- **Tokenization:** Converting raw text into tokens (words, punctuation, numbers). Essential for breaking text into analyzable units.
- **Lemmatization:** Converting each token to its base (dictionary) form.

- **Tokenization:** Converting raw text into tokens (words, punctuation, numbers). Essential for breaking text into analyzable units.
- **Lemmatization:** Converting each token to its base (dictionary) form.
- **Frequency Calculation:** Counting token/lemma occurrences in a text.

- **Tokenization:** Converting raw text into tokens (words, punctuation, numbers). Essential for breaking text into analyzable units.
- **Lemmatization:** Converting each token to its base (dictionary) form.
- **Frequency Calculation:** Counting token/lemma occurrences in a text.
- **Concordance:** Displaying each occurrence of a word/phrase with its surrounding context.

- **Tokenization:** Converting raw text into tokens (words, punctuation, numbers). Essential for breaking text into analyzable units.
- **Lemmatization:** Converting each token to its base (dictionary) form.
- **Frequency Calculation:** Counting token/lemma occurrences in a text.
- **Concordance:** Displaying each occurrence of a word/phrase with its surrounding context.

Now, let's get our hands dirty!