

APSC-5984 Lab 4: File system

Due: 2023-02-13 (Monday) 23:59:59

- APSC-5984 Lab 4: File system
 - 0. Overview
 - 1 Working directly (WD)
 - 1.1 `os.getcwd()`
 - 1.2 `os.listdir()`
 - 1.3 `os.chdir()`
 - 2. Path
 - 2.1 Absolute path vs. relative path
 - 2.2 Navigation using `os.path.join()`
 - 3. Interacting with files
 - 3.1 Create a file (mode `w`)
 - 3.2 Append to a file (mode `a`)
 - 3.3 Read a file (mode `r`)
 - 4. String processing
 - 4.1 String slicing
 - 4.2 String split
 - 4.3 String replace
 - 4.4 Regular expression (RE)

0. Overview

In this lab, you will learn basic Python commands to navigate the file system and manipulate files using the `os` library. Coupling the basic knowledge of interacting with the file system, you will be able to read and write data from/to files through Python.

1 Working directly (WD)

Here are the common Python `os` methods to interact with the file system:

- `os.getcwd()`: get the current WD
- `os.listdir()`: list the content of a directory (default: WD)
- `os.chdir()`: change the WD

Before you can access to the methods, remember to import the `os` library:

```
import os
```

1.1 `os.getcwd()`

To know where your current location, which is formally known as working directory (WD), is in the file system, you can use `os.getcwd()` command to print the path of the WD. `getcwd()` is an abbreviation of "get current working directory".

```
os.getcwd()  
# output: '/home/niche'
```

1.2 `os.listdir()`

You can use `os.listdir()` to list all the files and folders in the current WD. It is commonly used when you want to know your next step (e.g., which file to open) while you are exploring the file system.

```
os.listdir()  
# output: ['file_1.txt', 'file_2.txt', 'folder_1', 'folder_2']
```

1.3 `os.chdir()`

After knowing the current accessible folders, you can use `os.chdir()` to set the WD to a new folder. `chdir()` is an abbreviation of "change directory".

```
os.chdir('/home/niche/folder_1')  
os.getcwd()  
# output: '/home/niche/folder_1'
```

2. Path

A path is a string that represents the location of a file or folder in the file system. An example is the WD path you obtained from the `os.getcwd()`. There are several ways to interact with a path in Python:

- Absolute path: starts with the root directory `/`.
- Relative path: starts with the current WD.
- Tilde sign `~`: an alias for the home directory of the current user.
- One-dot sign `.`: an alias for the current WD.
- Double-dot sign `..`: an alias for the parent directory of the current WD.
- `os.path.join()`: join multiple paths together.

2.1 Absolute path vs. relative path

So far, the paths we printed were all absolute paths, where they starts with the root directory `/`. However, you can also use relative path to navigate the file system. For example, if you are in the folder `/home/niche`, you can use `os.chdir('folder_1')` to go to the subfolder `/home/niche/folder_1`. Here `folder_1` is the relative path to the current WD `/home/niche`.

```
os.getcwd()
# output: '/home/niche'
os.chdir('folder_1')
os.getcwd()
# output: '/home/niche/folder_1'
```

Similarly, if you are in the subfolder `/home/niche/folder_1`, you can use a double-dot symbol `..` to refer to its parent folder `/home/niche`. For example, you can use `os.chdir('..')` to go back to the parent folder `/home/niche`.

```
os.getcwd()
# output: '/home/niche/folder_1'
os.chdir('..')
os.getcwd()
# output: '/home/niche'
```

2.2 Navigation using `os.path.join()`

The `os.path.join()` method is a convenient way to join multiple paths together. It is recommended to use over explicitly typing the path in a string is because it is OS-agnostic. For example, if you are using Windows, the path separator is `\` instead of `/`. Using `os.path.join()` will automatically adjust the path separator based on your OS.

```
os.path.join('home', 'niche', 'folder_1')
# output: 'home/niche/folder_1' if you are using Linux
# output: 'home\\niche\\folder_1' if you are using Windows
```

3. Interacting with files

Now that you know how to navigate the file system and basic knowledge of paths, we can start to interact with files in Python!

3.1 Create a file (mode `w`)

To create a file, we can use the `open()` method. The `open()` method takes two arguments: the path to the file and the mode. The path just works as the way we learned in the previous section, it can be either absolute or relative path. The mode is a string that specifies how you want to interact with the file. The most common modes are:

- `r`: read-only mode
- `w`: write-only mode
- `a`: append-only mode

We will go through each mode in the following section. Let's start with the `w` mode. The `w` mode will create a new file if the file does not exist, or overwrite the existing file if the file already exists. We can create a file

named `file.txt` in the current WD using the following code:

```
f = open('file.txt', 'w')
f.close()
```

We always need to use `close()` to release the file resource after we are done with the file. Otherwise, the file will be locked and you will not be able to access it. Let's try writing something in this file:

```
f = open('file.txt', 'w')
f.write('Hello world!\n')
f.write('This ia the second line\n')
f.write('This ia the third line\n')
f.close()
```

It is noteworthy that the `write()` method will not automatically add an end of line (EOL) character (`\n`) at the end of the line. Therefore, we need to add it manually.

If you open the file `file.txt` in a text editor, you should see the following content:

```
Hello world!
This ia the second line
This ia the third line
```

3.2 Append to a file (mode `a`)

The `a` mode is similar to the `w` mode, except that it will append the content to the end of the file instead of overwriting the existing content. Let's try appending something to the file `file.txt`:

```
f = open('file.txt', 'a')
f.write('This ia the fourth line\n')
f.close()
```

Check the file again to see if the content is appended as expected.

3.3 Read a file (mode `r`)

The `r` mode is used to read the content of a file. There are two ways to read the content of a file:

- `read()` : read the entire content of the file.
- `readlines()`: store each line of the file as an element in a list.

Let's take a look of the first method `read()`:

```
f = open('file.txt', 'r')
content = f.read()
print(content)
f.close()
# output: 'Hello world!\nThis ia the second line\nThis ia the third
line\n'
```

You might notice that the entire file content, including the EOL characters `\n`, is stored in the variable `content`. This approach works fine when ones want to access the raw information, but it is not convenient when we want to process the content line by line. Hence, we can use the `readlines()` method:

```
f = open('file.txt', 'r')
lines = f.readlines()
print(lines)
f.close()
# output: ['Hello world!\n', 'This ia the second line\n', 'This ia the
third line\n']
```

The file contents are split into a list of lines automatically. You can use `len()` method to count how many lines are there in the file:

```
len(lines)
# output: 3
```

4. String processing

String processing is an essential skill to extract meaningful information from a text-based file. In this section, we will go through some common string processing methods.

4.1 String slicing

We have introduced the data type `string` and the related operators the previous lab. In this section, you will apply the operation to a list of strings. First, let's create an example list:

```
file_list = [
    "2022/09/08_trialA_trt1.txt",
    "2022/09/15_trialA_trt2.txt",
    "2022/09/28_trialB_trt1.txt",
    "2022/09/30_trialB_trt2.txt",
    "2022/10/01_trialC_trt1.txt",
    "2022/10/08_trialC_trt2.txt"]
```

In this example, you are provided with a list of file names. You might soon notice that the file names are coded in the following format: `YYYY/MM/DD_trialX_trtY.txt`. In this practice, we want to extract the

date, trial, and treatment information from the provided list.

We know that the file extension names are always occupied the last three characters of the file. We can use the indexing operation we learned earlier to extract the extension name:

```
file_list[0]
# output: '2022/09/08_trialA_trt1.txt'
file_list[0][-3:]
# output: 'txt'
```

The slicing index `-3:` means we want every character from the third last character to the end of the string. To generate a list of the extension names, you can use a `for` loop:

```
extension_list = []
for file in file_list:
    extension_list.append(file[-3:])
print(extension_list)
# output: ['txt', 'txt', 'txt', 'txt', 'txt', 'txt']
```

In this example, the `file` variable is assigned to each element (file name) in the list `file_list`. Then we use the `append()` method to add the extracted information (extension name) to the target list `extension_list`.

4.2 String split

Now, we want to break each string to extract other information such as trials and treatments. We noticed that they are separated by the underscore `_`. We can use the `split()` method that takes the separator (i.e., `_`) as the argument to split the string:

```
filename = file_list[0]
filename.split('_')
# output: ['2022/09/08', 'trialA', 'trt1.txt']
```

We can apply the same logic to the first element of the string to extract date information:

```
filename = file_list[0]
elements = filename.split('_')
date = elements[0]
print(date)
# output: '2022/09/08'
yyyymmdd = date.split("/")
print(yyyymmdd)
# output: ['2022', '09', '08']
```

4.3 String replace

From the split result, we can see that the treatment information is not in the format we want. Each substring `trt1.txt` still tails with the file extension name `.txt`.

```
filename = file_list[0]
elements = filename.split('_')
trt = elements[2]
print(trt)
# output: 'trt1.txt'
```

We have two strategies to remove the file extension name. The first one is to use a slicing operation:

```
new_trt = trt[:-4]
print(new_trt)
# output: 'trt1'
```

Or, we can use the `replace()` method to replace the file extension name with an empty string:

```
new_trt = trt.replace('.txt', '')
print(new_trt)
# output: 'trt1'
```

4.4 Regular expression (RE)

If the string pattern is more complicated than position-based slicing or replacement, you want to use regular expression to implement a more flexible pattern recognition. We will introduce the `re` module that provides a set of methods to implement regular expression.

As always, we need to import the `re` module first:

```
import re
```

The `re` module provides a set of methods to implement regular expression. In this lab, we will only use `re.findall()` method, which returns a list of all the non-overlapping matches in the string. The syntax of the method is:

```
re.findall(pattern, string)
```

The `pattern` is the regular expression pattern, and the `string` is the string to be searched. So, what is a regular expression pattern? Here is a table of some common regular expression patterns:

Pattern	Description
.	Matches any character except newline.
^	Matches the start of the string.
\$	Matches the end of the string or just before the newline at the end of the string.
*	Matches 0 or more (greedy) repetitions of the preceding RE.
+	Matches 1 or more (greedy) repetitions of the preceding RE.
?	Matches 0 or 1 (greedy) of the preceding RE.
\w	Matches any alphanumeric character including the underscore.
\W	Matches any non-alphanumeric character.
\d	Matches any numeric digit.
\D	Matches any non-numeric digit.
\s	Matches any whitespace character.
\S	Matches any non-whitespace character.
[]	Matches any character in the brackets.
{m}	Matches exactly m copies of the previous RE.
{m,n}	Matches from m to n (inclusive) copies of the previous RE.

It is toally normal to get confused by the table above when you first see it. Don't worry, we will go through some examples to understand how to use regular expression.

Here is an example:

```
filelines = [
    "File name: 2022/09/08_trialA_trt1.txt",
    "treatment A: (1, 3) and (2, 4)",
    "treatment B: (5, 1, 8)",
    "treatment C: (2, 4), (1, 9, 8)"]
```

The data `filelines` was intentionally made to be a little bit messy. We want to extract the treatment information in each pair of parenthesis.

```
trts = []
for line in filelines:
    value = re.findall(r"\([\d\s,]*\)", line)
    trts.append(value)
print(trts)
# output: [[], ['(1, 3)', '(2, 4)'], ['(5, 1, 8)'], ['(2, 4)', '(1, 9, 8)']]
```

Let's break down the RE pattern:

- `\(` and `\)`: matches a left parenthesis and a right parenthesis, respectively. The back slash `\` is used to escape the special character `(` and `)`.
- `[\d\s,]`: Since the treatment information is a list of numbers separated by commas, we can use the square brackets `[]` to match any character in the brackets. The pattern `[\d\s,]` means:
 - `\d`: matches any digit.
 - `\s`: matches any whitespace.
 - `,`: matches a comma.
- `*`: Since there is no specific number of digits, we can use the `*` to match 0 or more repetitions of the preceding RE.