

Information Retrieval

COMP 479 Project 1
Text preprocessing with NLTK DEMO

Lin Ling (40153877)

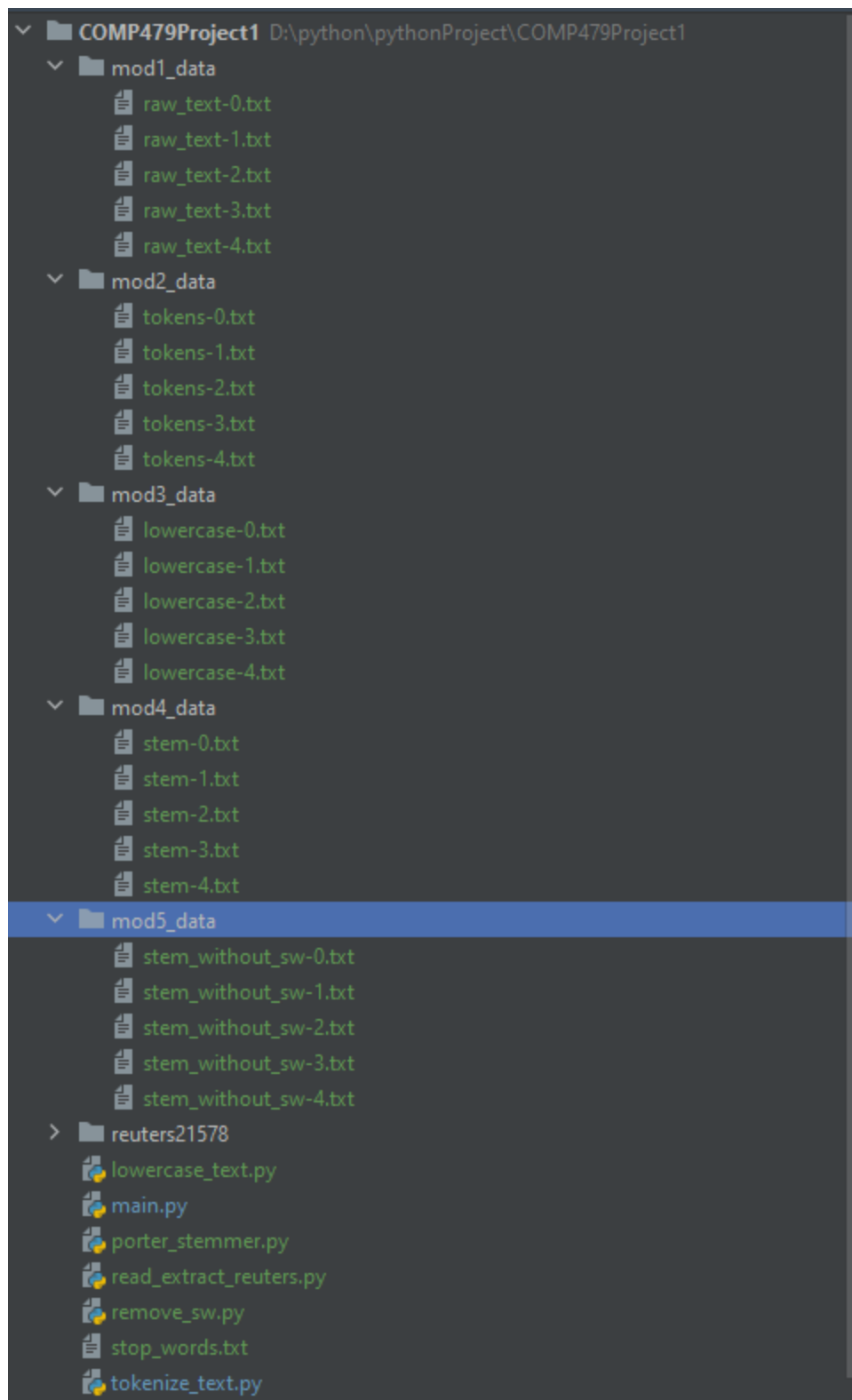
A report submitted in partial fulfilment of the requirements of Comp479.

Concordia University

Source Code: main.py

```
test.py x main.py x
1 import read_extract_reuters
2     import tokenize_text
3     import lowercase_text
4     import porter_stemmer
5 import remove_sw
6
7 TOP_FILES_NUM = 5
8 ORIGINAL_FILE_PATH = 'reuters21578'
9 ORIGINAL_FILE_TYPE = '.sgm'
10 MODULE_FILE_TYPE = '.txt'
11 MODULE_ONE_FILE_PATH = 'mod1_data'
12 MODULE_TWO_FILE_PATH = 'mod2_data'
13 MODULE_THREE_FILE_PATH = 'mod3_data'
14 MODULE_FOUR_FILE_PATH = 'mod4_data'
15 MODULE_FIVE_FILE_PATH = 'mod5_data'
16
17
18 # First Module: read the Reuter's collection and extract the raw text of each article from the corpus
19 print("Starting read the Reuter's collection ...")
20 read_extract_reuters.output_file(ORIGINAL_FILE_PATH, TOP_FILES_NUM, ORIGINAL_FILE_TYPE, MODULE_ONE_FILE_PATH)
21 print("Output the handled files in mod1_data directory.")
22
23 # Second Module: tokenize
24 print("Starting read the raw text files ...")
25 tokenize_text.output_file(MODULE_ONE_FILE_PATH, TOP_FILES_NUM, MODULE_FILE_TYPE, MODULE_TWO_FILE_PATH)
26 print("Output the handled files in mod2_data directory.")
27
28 # Third Module: make all text lowercase
29 print("Starting read the tokens files ...")
30 lowercase_text.output_file(MODULE_TWO_FILE_PATH, TOP_FILES_NUM, MODULE_FILE_TYPE, MODULE_THREE_FILE_PATH)
31 print("Output the handled files in mod3_data directory.")
32
33 # Fourth Module: apply Porter stemmer
34 print("Starting read the lowercase files ...")
35 porter_stemmer.output_file(MODULE_THREE_FILE_PATH, TOP_FILES_NUM, MODULE_FILE_TYPE, MODULE_FOUR_FILE_PATH)
36 print("Output the handled files in mod4_data directory.")
37
38 # Fifth Module: remove stop words
39 print("Starting read the stem files ...")
40 with open('stop_words.txt', 'r') as f:
41     stop_words = f.read().splitlines()
42 remove_sw.output_file(MODULE_FOUR_FILE_PATH, TOP_FILES_NUM, MODULE_FILE_TYPE, MODULE_FIVE_FILE_PATH, stop_words)
43 print("Output the handled files in mod5_data directory.")
44
45
```

Import five modules in a main.py to test the functionality of each part of the pipeline. All the read and write information like the number of files and the directory and path of the input files and output files is put as parameter in output_file method.



After implementation, each module output the files to their certain directory to easy manually proofread the result of steps in the pipeline.

```
main.py x remove_sw.py x lowercase_text.py x porter_stemmer.py x
1 import os
2
3
4 def remove_stop_words(path, words):
5     with open(path, "r") as f:
6         lines = f.read().splitlines()
7         start_len = len(lines)
8         lines = [line for line in lines if not line in words]
9         end_len = len(lines)
10        print(f'read {path} successfully.')
11        print(f'remove {start_len - end_len} words')
12        return lines
13
14
15 def output_file(input_path, file_num, file_ends, output_path, words):
16     count = 0
17     while count < file_num:
18         file_list = os.listdir(input_path)
19         file_list = [file for file in file_list if file.endswith(file_ends)]
20         for file in file_list[:file_num]:
21             stem_list = remove_stop_words(input_path + "/" + file, words)
22             with open(output_path + '/' + 'stem_without_sw-' + str(count) + '.txt', "a") as f:
23                 f.writelines("%s\n" % stem for stem in stem_list)
24             print(f'write to {output_path} successfully.')
25             count = count + 1
26
```

Put logs in each function's certain position to check the data flow in each step and easy to debug.

```
Run: main x
D:\Users\janel\anaconda3\python.exe D:/python/pythonProject/COMP479Project1/main.py
Starting read the stem files ...
read mod4_data/stem-0.txt successfully.
remove 20345 words
write to mod5_data successfully.
read mod4_data/stem-1.txt successfully.
remove 18934 words
write to mod5_data successfully.
read mod4_data/stem-2.txt successfully.
remove 17773 words
write to mod5_data successfully.
read mod4_data/stem-3.txt successfully.
remove 20747 words
write to mod5_data successfully.
read mod4_data/stem-4.txt successfully.
remove 20900 words
write to mod5_data successfully.
Output the handled files in mod5_data directory.
```

Do the test cases by unit test for each module by each function

```
test.py x main.py x
1 import unittest
2
3 from read_extract_reuters import read_split_article, extract_article_text
4 from tokenize_text import tokenize
5 from lowercase_text import lowercase
6 from porter_stemmer import porter_stemmer
7 from remove_sw import remove_stop_words
8
9
10 class TestModuleOne(unittest.TestCase):
11     def test_split_article(self):
12         """
13         Test that it can separate the whole text into each article
14         """
15         path_0 = "reuters21578/reut2-000.sgm"
16         path_21 = "reuters21578/reut2-021.sgm"
17         result = read_split_article(path_0)
18         self.assertEqual(len(result), 1000)
19         result = read_split_article(path_21)
20         self.assertEqual(len(result), 578)
21
22     def test_extract_article(self):
23         """
24         Test that it can extract the body content for each article
25         """
26         article_2 = "<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="16319" NEWID="999">
27 <DATE> 3-MAR-1987 09:16:02.85</DATE>
28 <TOPICS><D>money-fx</D><D>interest</D></TOPICS>
29 <PLACES><D>uk</D></PLACES>
30 <PEOPLE></PEOPLE>
31 <ORGS></ORGS>
32 <EXCHANGES></EXCHANGES>
33 <COMPANIES></COMPANIES>
34 <UNKNOWN>
```

Test results:

```
D:\Users\janel\anaconda3\python.exe "C:\Users\janel\AppData\Local\JetBrains\PyCharm Community Edition 2022.1.2\plugins\python-ce\helpers\pycharm\_jb_pytest_runner.py" --path "D:/t
Testing started at 9:41 p.m. ...
Launching pytest with arguments D:/0.Information Retrieval/Project/P1/COMP479Project1/test.py --no-header --no-summary -q in D:/0.Information Retrieval/Project/P1/COMP479Project1

===== test session starts =====
collecting ... collected 6 items

test.py::TestModuleOne::test_extract_article
test.py::TestModuleOne::test_split_article
test.py::TestModuleTwo::test_tokenize PASSED [ 16%]PASSED [ 33%]read reuters21578/reut2-000.sgm successfully.
read reuters21578/reut2-021.sgm successfully.

test.py::TestModuleThree::test_lowercase PASSED [ 50%]read mod1_data/raw_text-0.txt successfully.

test.py::TestModuleFour::test_porter_stemmer PASSED [ 66%]read mod2_data/tokens-0.txt successfully.

test.py::TestModuleFive::test_remove_stop_words PASSED [ 83%]read mod3_data/lowercase-0.txt successfully.
PASSED [100%]read mod4_data/stem-0.txt successfully.
remove 3699 words

===== 6 passed in 4.39s =====
```