# Information Retrieval

COMP 479 Project 2
Text preprocessing with NLTK DEMO

Lin Ling (40153877)

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

Concordia University

Source Code: naïve_index.py

```python
import util

# Subproject I: naive indexer
# # Reads documents and outputs term_documentID pairs
# term_IDs, term_IDs_positions = util.read_from_path("data")
#
# # sorts the list of term_documentsID pairs
# sorted_term_IDs = util.sorted_pairs(term_IDs)
# sorted_term_IDs_pos = util.sorted_triple(term_IDs_positions)
#
# # write the sorted pairs to a file
# util.write_pairs2file(sorted_term_IDs)
#
# # write the sorted triples to a file
# util.write_triples2file(sorted_term_IDs_pos)
#
# read file to a list of pairs
sorted_term_IDs = util.read_file2pairs("F.txt")

# removes duplicates in pairs list
unique_term_IDs = util.remove_dup_pairs(sorted_term_IDs)

# creates inverted index, structure of return dictionary {term. (frequency, posting_list)}
dict_tokens = util.create_posting_list(unique_term_IDs)
# Subproject II: single term query processing
print("--------------------Query processing for single term--------------------------")
print(util.query_term("outlook", dict_tokens))
print(util.query_term("zone", dict_tokens))
print(util.query_term("week", dict_tokens))

print("----Query processing for single term: four one-word challenge queries for Project 2-----")
print(util.query_term("copper", dict_tokens))
print(util.query_term("Samjens", dict_tokens))
print(util.query_term("Carmark", dict_tokens))
print(util.query_term("Bundesbank", dict_tokens))
```

Execute the naïve_index.py will print out the query processing results of three sample queries by myself and four one-word challenge queries.

Note: For saving time, comment out the code from line 3 to line 17, as it will take a long time to read 21578 files in data folder and generate the list of term_ids and term_ids_positions to a file which already run and got the file F.txt and F_triples.txt

Source code: effect_compression.py

After implementation, it will generate three tables to show the changed rate by different data retaining in different lossy compression steps. In the meantime, it will print out the query

processing results of three sample queries by myself and four one-word challenge queries after lossy compression steps.



```
Run:    effect_compression ×

--------------------Effect of preprocessing for Reuters-21578 Dictionary--------------------
+--------------+-------+-----+-----+
|              | size  |  ▲  | cml |
+--------------+-------+-----+-----+
|  unfiltered  | 50690 |     |     |
|  no_numbers  | 48027 |  -5 |  -5 |
| case folding | 39062 | -18 | -23 |
|  30 stopw's  | 39032 |  0  | -23 |
| 150 stopw's  | 38912 |  0  | -23 |
|   stemming   | 28578 | -26 | -49 |
+--------------+-------+-----+-----+

------------------Effect of preprocessing for Reuters-21578 non-positional index--------------------
+--------------+---------+-----+-----+
|              |  size   |  ▲  | cml |
+--------------+---------+-----+-----+
|  unfiltered  | 1640419 |     |     |
|  no_numbers  | 1480347 |  -9 |  -9 |
| case folding | 1431270 |  -3 | -12 |
|  30 stopw's  | 1159019 | -19 | -31 |
| 150 stopw's  |  881748 | -38 | -50 |
|   stemming   |  841413 |  -4 | -54 |
+--------------+---------+-----+-----+

------------------Effect of preprocessing for Reuters-21578 positional index--------------------
+--------------+---------+-----+-----+
|              |  size   |  ▲  | cml |
+--------------+---------+-----+-----+
|  unfiltered  | 2717928 |     |     |
|  no_numbers  | 2498724 |  -8 |  -8 |
| case folding | 2498724 |  0  |  -8 |
|  30 stopw's  | 1632989 | -34 | -42 |
| 150 stopw's  | 1178620 | -52 | -60 |
|   stemming   | 1178620 |  0  | -60 |
+--------------+---------+-----+-----+
--------------------Query processing for single term After Compression--------------------------
(191, [3330, 6259, 11732, 12375, 15551, 16823, 20461, 4, 16, 23, 54, 59, 144, 179, 209, 237, 259, 317, 481, 1152, 1214, 1252, 1405, 14
(82, [2522, 9852, 14747, 1323, 1990, 3442, 5167, 5244, 5273, 12502, 19110, 19559, 1, 273, 626, 630, 2264, 2411, 4246, 4442, 5057, 5143
(564, [95, 4533, 10178, 15248, 15978, 12732, 1, 4, 16, 109, 138, 144, 200, 203, 209, 249, 273, 278, 335, 339, 465, 518, 599, 612, 685,
----Query processing for single term After Compression: four one-word challenge queries for Project 2-----
(120, [22, 793, 800, 816, 922, 1148, 1184, 1552, 1607, 2006, 2074, 2186, 2764, 2782, 2880, 3454, 3613, 3862, 4058, 4291, 4373, 4431, 4
(5, [17837, 17863, 18069, 18071, 19419])
(1, [19758])
(167, [278, 926, 942, 950, 1540, 1959, 1969, 1971, 2110, 2197, 2662, 2686, 2979, 3020, 3514, 4062, 4113, 4297, 4828, 4830, 4873, 5176,
```
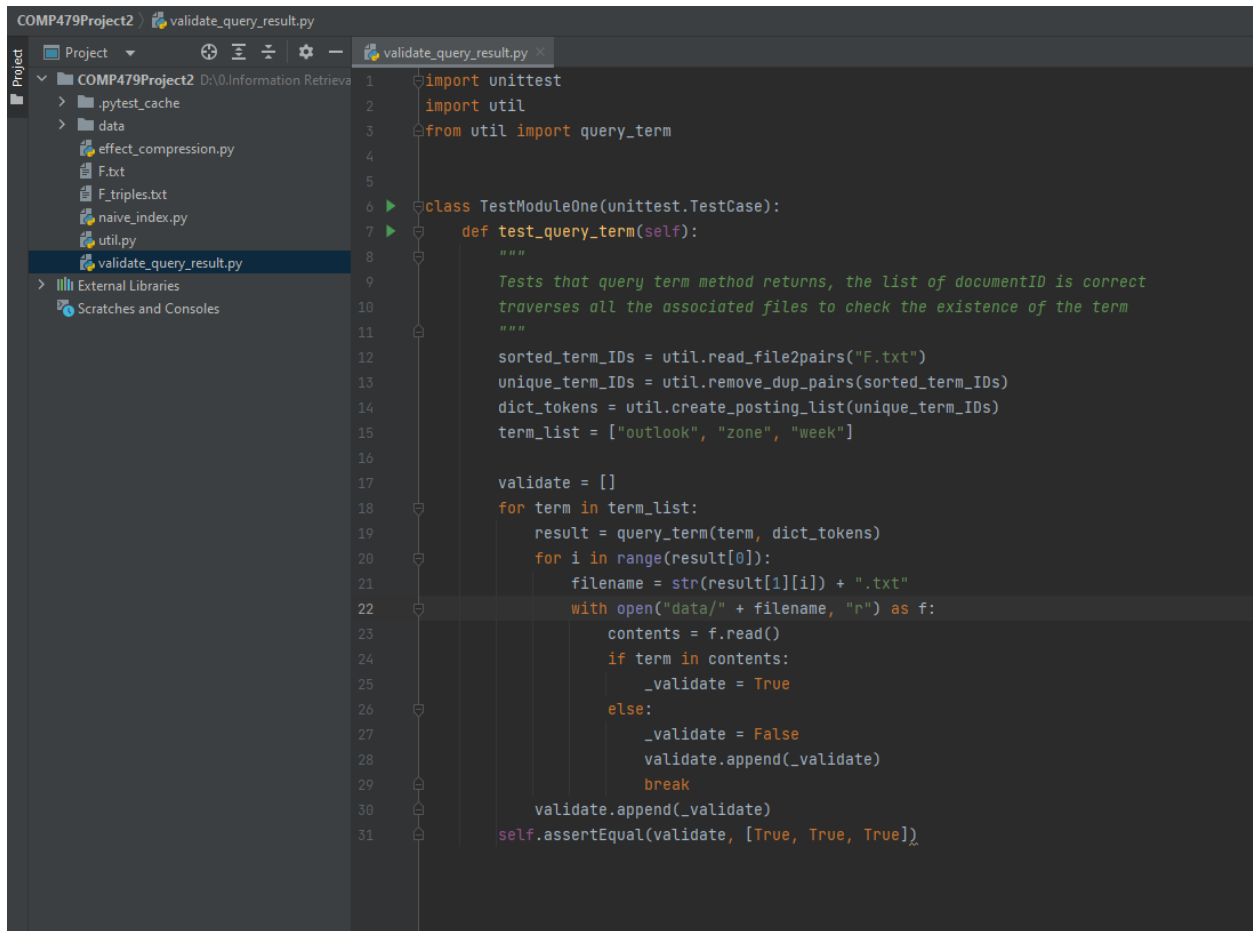
Source Code: util.py

This file contains all the function which used in naïve_index.py and effect_compression.py.

Source Code: Validate_query_result.py

Do the test case by unit test for validate query result



```python
import unittest
import util
from util import query_term


class TestModuleOne(unittest.TestCase):
    def test_query_term(self):
        """
        Tests that query term method returns, the list of documentID is correct
        traverses all the associated files to check the existence of the term
        """
        sorted_term_IDs = util.read_file2pairs("F.txt")
        unique_term_IDs = util.remove_dup_pairs(sorted_term_IDs)
        dict_tokens = util.create_posting_list(unique_term_IDs)
        term_list = ["outlook", "zone", "week"]

        validate = []
        for term in term_list:
            result = query_term(term, dict_tokens)
            for i in range(result[0]):
                filename = str(result[1][i]) + ".txt"
                with open("data/" + filename, "r") as f:
                    contents = f.read()
                    if term in contents:
                        _validate = True
                    else:
                        _validate = False
                        validate.append(_validate)
                        break
            validate.append(_validate)
        self.assertEqual(validate, [True, True, True])
```

Test results:

```
D:\Users\janel\anaconda3\python.exe "C:\Users\janel\AppData\Local\JetBrains\PyCharm Community Ed
Testing started at 10:46 p.m. ...
Launching pytest with arguments D:/0.Information Retrieval/Project/P2/COMP479Project2/validate_q

============================== test session starts ==============================
collecting ... collected 1 item

validate_query_result.py::TestModuleOne::test_query_term

============================== 1 passed in 6.33s ==============================

Process finished with exit code 0
PASSED            [100%]
```