**Cmpe 493 Introduction to Information Retrieval, Fall 2021 Assignment 1**
**A Simple Search System for Boolean Queries**
**Zuhal Didem Aytaç – 2018400045**

## 1. Data Preprocessing

I created a file preprocess.py. This code handles the initial processing of data. It reads the reutx-xxx.sgm files and the stopwords.txt file. It extracts the necessary content from the sgm files. Then with the helper functions, it performs case-fold, punctuation removal and stop-word removal operations on the extracted texts. It keeps the normalized, tokenized text using a dictionary, with new_id as the key.

This preprocess module is not run by itself, it is called by the indexing module.

## 2. The Inverted Index

The indexize module calls the preprocess module described in part 1. The preprocess module returns a dictionary (normalized) with new_id as key and list of tokens as value. The indexize module traverses that dictionary (normalized) and builds a new dictionary (index) with token as keys and list of new_ids as value. This dictionary (index) is a defaultdict, with the default value being set so that no new_ids are duplicated.
After that, the defaultdict (index) is copied into another dictionary (inverted_index) with values being sorted. Finally, the inverted_index dictionary is dumped to a json file (index.json) with token as key and sorted list of new_ids as value.
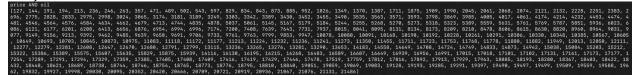
## 3. The Indexing Module

```
didemaytac@Zuhal-MacBook-Pro src % python3 indexize.py
```

## 4. Query Processing

You don't have to specify the query type. Type python3 query_processor.py, then type the query and press enter. The output will be shown at the command line.

```
didemaytac@Zuhal-MacBook-Pro src % python3 query_processor.py
```

Conjuction:

```
price AND oil
[127, 144, 191, 194, 213, 236, 246, 263, 357, 471, 489, 502, 543, 597, 829, 834, 843, 873, 885, 952, 1026, 1349, 1370, 1387, 1711, 1875, 1909, 1990, 2045, 2061, 2068, 2074, 2121, 2132, 2228, 2251, 2383, 2
696, 2775, 2828, 2833, 2975, 2998, 3024, 3065, 3174, 3181, 3189, 3249, 3303, 3342, 3389, 3430, 3452, 3455, 3490, 3535, 3563, 3571, 3593, 3798, 3869, 3985, 4005, 4017, 4061, 4174, 4214, 4232, 4453, 4474, 4
481, 4546, 4564, 4576, 4584, 4634, 4662, 4679, 4713, 4744, 4835, 4878, 5037, 5061, 5145, 5167, 5179, 5184, 5244, 5255, 5268, 5270, 5273, 5318, 5323, 5389, 5559, 5631, 5761, 5769, 5787, 5851, 5936, 6023, 6
086, 6121, 6177, 6201, 6208, 6413, 6656, 6876, 6954, 6994, 6996, 7174, 7200, 7408, 7639, 7643, 7731, 7937, 8015, 8041, 8095, 8131, 8134, 8173, 8209, 8210, 8478, 8606, 8615, 8630, 8820, 8960, 8964, 9031, 9
077, 9149, 9156, 9213, 9392, 9462, 9485, 9639, 9650, 9691, 9706, 9733, 9761, 9763, 9799, 9853, 9947, 10078, 10080, 10091, 10168, 10190, 10192, 10228, 10261, 10291, 10306, 10330, 10348, 10385, 10567, 10605
, 10649, 10693, 10703, 10845, 10873, 10975, 11083, 11118, 11172, 11177, 11213, 11224, 11232, 11236, 11241, 11273, 11350, 11455, 11711, 11723, 11753, 11768, 11778, 11880, 11882, 11949, 12013, 12050, 12111,
 12277, 12279, 12281, 12608, 12647, 12670, 12680, 12791, 12799, 13115, 13236, 13265, 13276, 13281, 13290, 13653, 14183, 14558, 14649, 14708, 14724, 14749, 14833, 14873, 14942, 15038, 15084, 15203, 15212,
15322, 15386, 15389, 15575, 15607, 15635, 15829, 15875, 15939, 16116, 16130, 16195, 16215, 16268, 16483, 16589, 16607, 16649, 16939, 16956, 16991, 17015, 17018, 17101, 17102, 17131, 17161, 17173, 17177, 1
7254, 17289, 17291, 17294, 17329, 17359, 17385, 17405, 17408, 17409, 17416, 17419, 17429, 17446, 17478, 17519, 17759, 17812, 17816, 17892, 17913, 17929, 17963, 18085, 18193, 18280, 18367, 18403, 18422, 18
432, 18448, 18621, 18689, 18738, 18744, 18746, 18754, 18765, 18773, 18776, 18795, 18810, 18840, 19051, 19059, 19069, 19083, 19128, 19193, 19285, 19291, 19397, 19490, 19497, 19499, 19509, 19559, 19588, 196
62, 19832, 19927, 19998, 20030, 20095, 20352, 20420, 20666, 20709, 20721, 20919, 20936, 21067, 21076, 21131, 21486]
```

## Disjunction:

petroleum OR oil OR gas
[2, 6, 8, 26, 68, 84, 91, 127, 137, 140, 144, 145, 156, 157, 176, 191, 194, 200, 211, 213, 235, 236, 237, 242, 246, 247, 248, 263, 273, 274, 277, 288, 298, 304, 313, 320, 332, 340, 349, 352, 353, 356, 357, 364, 368, 370, 372, 391, 450, 459, 471, 489, 502, 507, 542, 543, 544, 570, 593, 597, 613, 622, 666, 668, 697, 704, 708, 739, 791, 799, 829, 834, 835, 837, 843, 855, 862, 873, 885, 888, 896, 915, 918, 92
...
1525, 21536, 21541, 21561, 21568]

## Conjunction and Negation:

price AND oil NOT vegetable
[127, 144, 191, 194, 236, 246, 263, 357, 471, 489, 502, 543, 597, 829, 834, 843, 873, 885, 952, 1026, 1349, 1370, 1387, 1711, 1875, 1909, 1990, 2045, 2061, 2068, 2074, 2121, 2132, 2228, 2251, 2383, 2696, 2775, 2828, 2833, 2975, 2998, 3024, 3065, 3174, 3181, 3189, 3249, 3303, 3342, 3389, 3430, 3452, 3455, 3490, 3535, 3563, 3571, 3593, 3798, 3869, 3985, 4005, 4017, 4061, 4174, 4214, 4232, 4453, 4474, 4481,
...
0352, 20420, 20666, 20709, 20721, 20919, 20936, 21067, 21076, 21131, 21486]

## Disjunction and Negation:

petroleum OR oil NOT price
[2, 6, 8, 26, 68, 137, 140, 145, 156, 157, 200, 211, 235, 237, 242, 247, 248, 273, 274, 277, 288, 313, 320, 332, 340, 349, 352, 353, 356, 364, 368, 370, 391, 459, 542, 570, 593, 613, 666, 668, 697, 704, 7
08, 739, 791, 835, 837, 855, 888, 896, 915, 918, 927, 930, 939, 944, 945, 957, 963, 978, 988, 1004, 1024, 1046, 1084, 1098, 1110, 1127, 1140, 1150, 1211, 1215, 1297, 1301, 1306, 1316, 1335, 1343, 1379, 14
...
21492, 21501, 21502, 21506, 21510, 21519, 21525, 21541, 21561, 21568]