

## Assignment 1

Describe the steps you have performed for data preprocessing.

In this project we are assigned to implement a document retrieval system for simple boolean queries using the non-positional inverted indexing scheme. Reuters news is our dataset, and we have some stopwords to discard. First, I stored stopwords in a list for not to include these words while indexing texts. I have a dictionary named `index_dict`. I store newId's of each word in here. Then I started to traverse .sgm files in reuters21578 directory. I used BeautifulSoup library to ease .sgm file processing. Then, I store all articles in a file in `news_list` list. I iterate over each news and look for their titles and body's. If the article contains either of these, I combine the title and body and send it to the `normalize` function. In `normalize` function I made punctuation removal, case folding and finally stopwords removal. It also prevents duplicate entries in text. After these steps I added these words and their ids to dictionary and printed this dictionary to file.

Describe the data structures (hash, b-tree, linked list etc.) that you used for representing the inverted index (i.e., the dictionary and the postings lists).

I used list while storing posting lists of words. Since .sgm files read in ascending order of news ids it stores new ids in ascending order. While implementing AND operator I used pseudocode in Lec01. I traversed two list by comparing news ids and make necessary addition to result.

Provide a screenshot of running the indexing module of your system.

Since I am using BeautifulSoup library it must be installed on working environment.

```
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment1> pip install bs4
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
Requirement already satisfied: bs4 in c:\python310\lib\site-packages (0.0.1)
Requirement already satisfied: beautifulsoup4 in c:\python310\lib\site-packages (from bs4) (4.10.0)
Requirement already satisfied: soupsieve>1.2 in c:\python310\lib\site-packages (from beautifulsoup4->bs4) (2.3.1)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution -ip (c:\python310\lib\site-packages)
WARNING: Ignoring invalid distribution - (c:\python310\lib\site-packages)
WARNING: You are using pip version 21.2.3; however, version 21.3.1 is available.
You should consider upgrading via the 'C:\Python310\python.exe -m pip install --upgrade pip' command.
PS C:\Users\muhammed\Desktop\Ders\7.Dönem\Cmpe493\Assignment1> python .\preprocessing.py
```

## Muhammed Göktepe 2017400162

You can install it by typing “pip install bs4”. Then you can run my preprocessing code by typing “python preprocessing.py”.

Provide four screenshots of running your system for each of the four types of queries.  
price AND oil

```
PS C:\Users\muhammed\Desktop\Users\7.Dönem\Cnpe493\Assignment1> python .\queryprocessing.py
Please enter query:price AND oil
127, 144, 191, 194, 215, 236, 240, 243, 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, 3069, 3174, 3181, 3189, 3249, 3303, 3342, 3389, 3430, 3452, 3495, 3498, 3535, 3563, 3571, 3593, 3798, 3869, 3945, 4005, 4017, 4061, 4174, 4214, 4232, 4453, 4474, 4481, 4546, 4564, 4576, 4584, 4634, 4682, 4679, 4713, 474
7, 4835, 4878, 5037, 5041, 5145, 5149, 5179, 5184, 5246, 5255, 5268, 5270, 5319, 5323, 5389, 5559, 5631, 5741, 5767, 5787, 5851, 5935, 6023, 6084, 6121, 6177, 6201, 6280, 6413, 6856, 6876, 6954, 6994, 6998, 7174, 7280, 7488, 7
39, 7643, 7711, 7837, 8015, 8041, 8095, 8111, 8134, 8172, 8209, 8210, 8478, 8486, 8619, 8630, 8820, 8846, 8944, 9031, 9077, 9149, 9154, 9215, 9292, 9462, 9485, 9639, 9640, 9651, 9766, 9735, 9741, 9743, 9799, 9833, 9847, 10078, 1008
9, 10091, 10168, 10198, 10192, 10228, 10261, 10291, 10300, 10330, 10348, 10385, 10547, 10605, 10649, 10693, 10703, 10845, 10871, 10885, 10923, 10975, 11063, 11118, 11172, 11177, 11213, 11226, 11232, 11236, 11241, 11271, 11330, 11455, 11711, 11723
11753, 11768, 11778, 11880, 11882, 11949, 12813, 12958, 12111, 12277, 12279, 12281, 12608, 12647, 12678, 12680, 12791, 12799, 13115, 13236, 13265, 13276, 13281, 13290, 13453, 14183, 14558, 14649, 14708, 14724, 14749, 14833, 14873,
14942, 15038, 15084, 15203, 15212, 15322, 15352, 15385, 15399, 15575, 15607, 15635, 15829, 15875, 15939, 16116, 16130, 16195, 16215, 16268, 16483, 16589, 16607, 16649, 16939, 16956, 16991, 17015, 17018, 17101, 17102, 17311, 17161, 17173,
17177, 17254, 17289, 17291, 17294, 17329, 17359, 17385, 17405, 17408, 17409, 17416, 17419, 17429, 17446, 17478, 17519, 17579, 17812, 17816, 17892, 17913, 17929, 17963, 18085, 18193, 18280, 18367, 18403, 18422, 18432, 18448, 18621, 1
8689, 18738, 18744, 18746, 18754, 18765, 18773, 18778, 18795, 18810, 18840, 19051, 19059, 19069, 19083, 19128, 19193, 19285, 19291, 19397, 19490, 19497, 19499, 19509, 19559, 19588, 19662, 19832, 19927, 19998, 20030, 20095, 20352, 20
320, 20566, 20709, 20721, 20919, 20936, 21067, 21076, 21131, 21468]
```

petroleum OR oil OR gas

```
PS C:\Users\muhammed\Desktop\Users\7.Dönem\Cnpe493\Assignment1> python .\queryprocessing.py
Please enter query:petroleum OR oil OR gas
12, 6, 8, 26, 08, 84, 91, 127, 137, 140, 144, 149, 156, 159, 159, 170, 191, 194, 200, 211, 215, 235, 236, 237, 242, 246, 247, 249, 263, 273, 274, 277, 288, 298, 304, 313, 320, 332, 349, 349, 359, 352, 353, 355, 357, 364, 368, 370, 371, 391, 450, 45
9, 471, 489, 502, 509, 542, 543, 546, 570, 573, 597, 613, 622, 646, 668, 697, 704, 708, 739, 741, 759, 829, 836, 835, 837, 843, 855, 862, 873, 885, 888, 890, 919, 919, 927, 938, 939, 946, 945, 952, 957, 943, 978, 988, 1004, 1024, 1026, 10
46, 1064, 1084, 1098, 1112, 1127, 1140, 1150, 1211, 1212, 1297, 1301, 1380, 1316, 1330, 1343, 1349, 1370, 1379, 1387, 1406, 1463, 1440, 1501, 1520, 1521, 1550, 1552, 1556, 1558, 1616, 1619, 1680, 1688, 1690, 1681, 1688, 1692, 1696, 1703,
1709, 1711, 1723, 1751, 1756, 1780, 1799, 1824, 1825, 1837, 1856, 1868, 1874, 1875, 1878, 1891, 1900, 1909, 1957, 1948, 1959, 1964, 1980, 1990, 1999, 2002, 2007, 2045, 2046, 2061, 2068, 2074, 2121, 2132, 2159, 2173, 2175, 2187, 2218, 2213
, 2221, 2228, 2231, 2234, 2231, 2262, 2268, 2303, 2383, 2394, 2415, 2423, 2432, 2435, 2446, 2460, 2475, 2477, 2479, 2480, 2483, 2495, 2511, 2515, 2517, 2522, 2528, 2530, 2542, 2582, 2585, 2596, 2685, 2688, 2696, 2737, 2738, 2745, 27
58, 2767, 2775, 2780, 2789, 2811, 2818, 2820, 2822, 2828, 2833, 2834, 2835, 2838, 2897, 2923, 2925, 2927, 2930, 2957, 2970, 2973, 2975, 2998, 3003, 3015, 3017, 3019, 3024, 3048, 3065, 3080, 3115, 3145, 3146, 3169, 3174, 3181, 3189, 3204,
3260, 3210, 3244, 3252, 3264, 3269, 3303, 3310, 3332, 3338, 3342, 3344, 3354, 3364, 3372, 3389, 3411, 3430, 3434, 3452, 3455, 3462, 3466, 3468, 3496, 3505, 3507, 3509, 3535, 3548, 3550, 3554, 3571, 3592, 3593, 3594, 3597, 3609, 3615, 3657
, 3714, 3736, 3747, 3752, 3758, 3763, 3779, 3787, 3798, 3800, 3818, 3848, 3845, 3848, 3855, 3855, 3864, 3869, 3872, 3888, 3906, 3925, 3929, 3940, 3950, 3980, 3985, 3986, 3995, 4005, 4016, 4017, 4027, 4028, 4037, 4039, 4041, 4049, 4061, 40
07, 4089, 4125, 4126, 4129, 4134, 4138, 4152, 4171, 4174, 4189, 4201, 4209, 4214, 4232, 4246, 4249, 4290, 4305, 4315, 4316, 4325, 4333, 4337, 4338, 4340, 4353, 4359, 4359, 4384, 4425, 4429, 4453, 4464, 4467, 4474, 4479, 4484, 4491, 4503,
4507, 4510, 4522, 4530, 4531, 4540, 4547, 4553, 4554, 4566, 4569, 4570, 4577, 4584, 4589, 4590, 4593, 4600, 4601, 4604, 4609, 4634, 4650, 4661, 4662, 4679, 4681, 4687, 4702, 4713, 4730, 4734, 4742, 4744, 4753, 4785, 4834, 4835, 4848, 4867
, 4878, 4900, 4908, 4930, 4936, 4947, 4948, 4951, 4953, 4962, 4963, 4981, 4983, 5030, 5035, 5037, 5044, 5041, 5046, 5110, 5118, 5119, 5123, 5125, 5137, 5142, 5143, 5145, 5158, 5152, 5156, 5163, 5165, 5166, 5169, 5167, 5178, 5179, 51
84, 5187, 5193, 5203, 5210, 5218, 5230, 5238, 5244, 5250, 5255, 5268, 5270, 5273, 5274, 5281, 5282, 5295, 5315, 5318, 5323, 5330, 5342, 5370, 5389, 5429, 5440, 5449, 5472, 5485, 5506, 5538, 5541, 5542, 5544, 5552, 5553, 5558, 5559, 5561, 5578,
5630, 5631, 5655, 5656, 5669, 5675, 5683, 5684, 5688, 5692, 5700, 5710, 5712, 5739, 5751, 5755, 5761, 5769, 5776, 5787, 5789, 5791, 5793, 5796, 5828, 5830, 5848, 5851, 5852, 5866, 5879, 5887, 5919, 5920, 5936, 5949, 5953, 5965, 5965, 6001
, 6023, 6037, 6049, 6052, 6054, 6060, 6080, 6087, 6098, 6111, 6119, 6121, 6125, 6133, 6145, 6159, 6162, 6163, 6166, 6169, 6177, 6184, 6193, 6201, 6208, 6219, 6225, 6258, 6264, 6271, 6286, 6301, 6380, 6388, 6317, 6322, 6342, 6344, 6348, 63
54, 6371, 6484, 6412, 6413, 6421, 6425, 6432, 6435, 6535, 6560, 6562, 6573, 6578, 6596, 6606, 6638, 6650, 6652, 6656, 6660, 6670, 6685, 6688, 6692, 6708, 6712, 6722, 6740, 6742, 6746, 6760, 6836, 6856, 6869, 6873, 6879, 6879, 6888, 6893,
6965, 6927, 6923, 6954, 6967, 6974, 6996, 7003, 7050, 7044, 7047, 7097, 7117, 7135, 7150, 7152, 7154, 7174, 7190, 7200, 7287, 7291, 7307, 7312, 7317, 7355, 7356, 7367, 7408, 7416, 7423, 7441, 7462, 7469, 7508, 7515, 7518, 7529, 7534, 7547
, 7548, 7589, 7608, 7611, 7615, 7618, 7625, 7639, 7642, 7643, 7653, 7671, 7702, 7721, 7740, 7742, 7745, 7750, 7831, 7854, 7855, 7883, 7894, 7894, 7901, 7937, 7996, 8005, 8014, 8019, 8021, 8033, 8039, 8041, 8042, 8060, 8081, 8085, 8082, 80
99, 8088, 8089, 8095, 8100, 8108, 8109, 8117, 8119, 8129, 8131, 8134, 8149, 8156, 8159, 8160, 8167, 8173, 8188, 8190, 8209, 8210, 8237, 8238, 8239, 8288, 8331, 8335, 8339, 8402, 8405, 8414, 8421, 8423, 8440, 8478, 8493, 8502, 8516, 8530,
8547, 8570, 8594, 8598, 8600, 8610, 8615, 8623, 8630, 8633, 8672, 8675, 8688, 8703, 8747, 8755, 8764, 8765, 8788, 8812, 8815, 8828, 8824, 8835, 8854, 8856, 8877, 8882, 8884, 8894, 8920, 8929, 8959, 8960, 8964, 8971, 8988, 8995
, 9031, 9045, 9077, 9098, 9100, 9112, 9127, 9149, 9155, 9156, 9180, 9185, 9193, 9204, 9206, 9208, 9213, 9253, 9256, 9260, 9275, 9279, 9293, 9328, 9352, 9354, 9370, 9381, 9392, 9402, 9436, 9445, 9462, 9475, 9478, 9479, 9483, 9485, 9527, 95
44, 9545, 9547, 9549, 9550, 9573, 9583, 9601, 9610, 9614, 9634, 9639, 9645, 9650, 9674, 9691, 9692, 9708, 9718, 9722, 9731, 9733, 9734, 9736, 9737, 9742, 9756, 9761, 9763, 9769, 9770, 9778, 9799, 9801, 9821, 9845, 9849, 9859, 9853, 9872, 9908,
9913, 9947, 9952, 10011, 10045, 10078, 10080, 10091, 10100, 10123, 10135, 10137, 10146, 10147, 10150, 10168, 10170, 10175, 10190, 10192, 10200, 10226, 10228, 10261, 10268, 10272, 10275, 10284, 10291, 10292, 10300, 10306, 10317, 103
30, 10341, 10348, 10345, 10371, 10373, 10375, 10385, 10386, 10395, 10405, 10408, 10440, 10442, 10521, 10539, 10547, 10551, 10588, 10605, 10620, 10621, 10624, 10627, 10628, 10639, 10632, 10641, 10649, 10670, 10693, 10703, 10719, 10720, 107
59, 10783, 10792, 10797, 10812, 10814, 10831, 10845, 10848, 10847, 10849, 10862, 10873, 10898, 10910, 10927, 10944, 10947, 10952, 10955, 10981, 11000, 11002, 11007, 11025, 11044, 11053, 11056, 11062, 11070, 11083, 11089, 11093, 11109, 11
61, 11118, 11136, 11145, 11149, 11168, 11171, 11172, 11177, 11213, 11224, 11227, 11231, 11232, 11236, 11237, 11241, 11249, 11270, 11273, 11275, 11281, 11305, 11325, 11350, 11351, 11386, 11393, 11403, 11406, 11421, 11444, 11455, 11465, 1114
66, 11487, 11542, 11526, 11522, 11533, 11559, 11588, 11600, 11632, 11635, 11636, 11639, 11671, 11674, 11682, 11697, 11702, 11705, 11710, 11711, 11717, 11723, 11724, 11728, 11731, 11732, 11752, 11753, 11760, 11768, 11778, 11781, 11782, 118
12, 11822, 11872, 11888, 11882, 11886, 11892, 11898, 11906, 11908, 11943, 11949, 11950, 11953, 12013, 12029, 12031, 12090, 12060, 12111, 12118, 12119, 12122, 12123, 12149, 12151, 12154, 12165, 12173, 12181, 12193, 12209, 12224, 12241, 1221
77, 12279, 12281, 12280, 12289, 12300, 12312, 12313, 12320, 12364, 12365, 12361, 12369, 12378, 12384, 12454, 12563, 12533, 12521, 12531, 12533, 12536, 12578, 12595, 12606, 12608, 12608, 12609, 12611, 12636, 12647, 12670, 12680, 126
86, 12691, 12708, 12709, 12717, 12728, 12739, 12746, 12755, 12775, 12778, 12794, 12799, 12803, 12818, 12822, 12831, 12845, 12855, 12859, 12953, 12975, 12977, 12986, 12988, 13025, 13035, 13051, 13056, 13060, 13074, 13076, 13080, 130
85, 13090, 13102, 13115, 13127, 13138, 13162, 13179, 13181, 13184, 13208, 13206, 13236, 13245, 13347, 13350, 13385, 13401, 13420, 13422, 13501, 13517, 13539, 13542, 13576, 13601, 136
13, 13654, 13653, 13702, 13713, 13749, 13753, 13833, 13847, 13858, 13861, 13907, 13915, 13949, 13951, 14001, 14102, 14107, 14183, 14190, 14211, 14230, 14248, 14279, 14299, 14395, 14419, 14500, 14509, 14544, 14555, 14585, 14611, 146
16, 14635, 14646, 14667, 14679, 14680, 14681, 14690, 14698, 14700, 14708, 14716, 14726, 14732, 14734, 14742, 14749, 14755, 14760, 14779, 14822, 14832, 14833, 14840, 14853, 14863, 14873, 14894, 14913, 14922, 14929, 14934, 14942, 14972, 14980, 150
84, 15018, 15038, 15043, 15058, 15068, 15162, 15198, 15200, 15203, 15212, 15227, 15230, 15238, 15246, 15262, 15275, 15296, 15322, 15325, 15411, 15444, 15447, 15449, 15491, 15493, 15506, 15538, 15580, 15589, 15596, 15615, 1565
75, 15599, 15607, 15613, 15635, 15655, 15673, 15759, 15760, 15812, 15814, 15824, 15829, 15853, 15871, 15875, 15890, 15900, 15909, 15923, 15924, 15931, 15939, 15944, 15975, 16000, 16005, 16007, 16044, 16055, 16059, 16077, 16080, 16086, 160
87, 16093, 16112, 16113, 16130, 16155, 16164, 16166, 16169, 16176, 16195, 16212, 16214, 16226, 16228, 16239, 16249, 16252, 16268, 16270, 16300, 16316, 16323, 16333, 16349, 16349, 16358, 16416, 16429, 16438, 16463, 16475, 164
83, 16499, 16508, 16532, 16577, 16589, 16593, 16597, 16604, 16607, 16613, 16636, 16640, 16644, 16649, 16651, 16658, 16665, 16680, 16702, 16710, 16713, 16723, 16731, 16733, 16739, 16742, 16762, 16773, 16789, 16807, 16830, 16833, 16894, 169
16, 16918, 16935, 16945, 16948, 16956, 16961, 16966, 16968, 16979, 16991, 16997, 16998, 17000, 17003, 17012, 17015, 17018, 17028, 17046, 17048, 17054, 17079, 17080, 17082, 17093, 17096, 17100, 17101, 17110, 17112, 17120, 171
26, 17131, 17138, 17139, 17161, 17170, 17173, 17177, 17180, 17185, 17186, 17190, 17199, 17208, 17225, 17236, 17245, 17246, 17249, 17291, 17294, 17296, 17327, 17329, 17333, 17353, 17354, 17355, 17387, 17395, 17401, 17408, 17369, 17349, 173
72, 17385, 17408, 17408, 17409, 17416, 17415, 17416, 17429, 17422, 17427, 17429, 17430, 17431, 17433, 17441, 17446, 17462, 17468, 17463, 17469, 17473, 17477, 17478, 17492, 17510, 17515, 17516, 17519, 17546, 17567, 17584, 17605, 17611, 176
18, 17653, 17678, 17682, 17693, 17708, 17735, 17759, 17771, 17780, 17797, 17812, 17813, 17816, 17848, 17874, 17878, 17883, 17888, 17888, 17892, 17894, 17911, 17913, 17924, 17928, 17929, 17929, 17941, 17968, 17963, 17968, 17968, 179
83, 18060, 18089, 18108, 18137, 18139, 18146, 18150, 18180, 18193, 18201, 18213, 18228, 18231, 18234, 18253, 18271, 18280, 18308, 18325, 18338, 18359, 18372, 18382, 18383, 18387, 18398, 18400, 18409, 18492, 18498, 18701, 18704, 18705, 18708, 18707, 18723, 18728, 187
19, 18422, 18431, 18432, 18444, 18447, 18448, 18471, 18480, 18493, 18504, 18521, 18523, 18547, 18573, 18579, 18621, 18624, 18637, 18651, 18655, 18657, 18678, 18680, 18
```

petroleum OR oil NOT price