

Report

Q1:

- I used re to remove symbols and other non-useful characters from the text data.
- I used nltk library to then remove punctuation, stop words and blank spaces from the text.
- 5 examples from different files.

```
File file270.txt content is:
let low price fool incredible device free mixing software years since tracked anything back everything solid state analogue scarlett solo moderate computer large monitor stress large m

File file215.txt content is:
great taylor 150e 12string guitar issue wand came tear microfiber duster

File file701.txt content is:
absolute junk doesnt worth amazon shipment price dont buy circumstances fire hazard 3 4 received dead arrival even missing parts rattle inside one working also something rattling insid

File file693.txt content is:
works perfectly ukuleles guitars simple use instant perfect tuner

File file302.txt content is:
case awesome molded interior fits fender jazz bass like glove feel secured times second case ive owned quality superb balances much better hand carrying last case bit front heavy also
```

Q2:

- I used the default dictionary to store the inverted indexes.
- Then I used a for loop to store all the document ids with corresponding tokens in the dictionary.

```
way: [1, 108, 13, 139, 144, 152, 157, 194, 214, 228, 249, 269, 279, 293, 305, 326, 349, 356, 373, 383, 384, 386, 406, 413, 428, 43, 438, 469, 483, 487, 49, 512, 537, 546, 549,
great: [1, 100, 101, 102, 103, 104, 105, 108, 109, 115, 119, 120, 121, 122, 124, 127, 129, 13, 131, 132, 134, 137, 139, 140, 147, 148, 152, 156, 158, 167, 176, 183, 187, 188, 1
go: [1, 100, 103, 106, 117, 136, 139, 144, 150, 172, 19, 191, 196, 212, 230, 24, 248, 254, 257, 332, 335, 353, 363, 366, 382, 388, 407, 410, 413, 414, 42, 423, 429, 454, 464, 4
floating: [1, 245, 430]
loving: [1, 254, 391, 723]
good: [1, 103, 106, 110, 111, 115, 118, 13, 137, 141, 143, 154, 155, 157, 159, 16, 160, 162, 163, 164, 166, 172, 174, 175, 176, 179, 18, 189, 19, 2, 204, 207, 210, 217, 220, 23
strat: [1, 149, 163, 197, 241, 245, 25, 253, 345, 353, 380, 396, 400, 422, 440, 455, 457, 469, 519, 529, 559, 565, 579, 611, 626, 650, 652, 691, 801, 838, 853, 90, 940, 978, 99
bridge: [1, 108, 117, 139, 168, 244, 249, 257, 283, 328, 353, 361, 454, 464, 484, 503, 541, 579, 598, 61, 620, 700, 725, 748, 756, 770, 787, 801, 806, 881, 916, 927, 998]
vintage: [1, 150, 197, 278, 422, 439, 494, 51, 597, 638, 674, 725, 737, 827, 847, 895, 907, 936]
tension: [1, 116, 143, 507, 514, 751, 839, 980]
springs: [1, 272, 469, 806, 937]
stability: [1, 115, 382, 521, 759]
want: [1, 102, 114, 121, 129, 143, 191, 213, 216, 228, 23, 232, 249, 293, 299, 307, 320, 325, 33, 343, 353, 359, 366, 39, 390, 396, 399, 419, 429, 43, 439, 490, 5, 519, 529, 54
product: [10, 105, 115, 121, 134, 142, 152, 167, 168, 171, 178, 180, 186, 19, 196, 208, 217, 225, 228, 229, 241, 251, 26, 269, 282, 298, 30, 326, 337, 35, 358, 363, 372, 374, 3
frame: [10, 352, 466, 684, 849]
bend: [10, 212, 70, 833, 994]
making: [10, 111, 143, 158, 172, 212, 249, 336, 338, 425, 444, 5, 744, 769, 871]
stand: [10, 120, 122, 135, 150, 153, 172, 18, 185, 208, 210, 234, 267, 269, 283, 286, 3, 303, 308, 326, 366, 374, 382, 408, 42, 425, 44, 442, 458, 459, 462, 465, 466, 47, 475,
angle: [10, 103, 301, 305, 338, 523, 579, 648, 663, 678, 720, 765, 772, 778, 79, 808, 882]
weird: [10, 499, 690, 721, 794, 886, 923]
photos: [10, 117, 126, 222, 290, 294, 319, 449, 460, 489, 494, 525, 55, 627, 694, 726, 790]
support: [10, 267, 270, 282, 42, 499, 646, 649, 778, 875, 937, 993]
vertical: [10]
becoming: [10, 254, 491]
assembled: [10, 458, 728]
```

- This is an example of the view of invert index

- This is an example query that I ran.

```
Number of documents retrieved for query 1 using positional index: 2
```

- Names of documents retrieved for query 1 using positional index: ['249.txt', '361.txt']