

```
#!/usr/bin/env python3

from collections import defaultdict

name = ""
inv_index = defaultdict(list)
postings = defaultdict(set)
stop_words = {'a', 'an', 'and', 'are', 'as', 'at', 'be', 'by', 'for', 'from',
              'has', 'he', 'in', 'is', 'it', 'its', 'of', 'on', 'that', 'the',
              'to', 'was', 'were', 'will', 'with'}

def normalize(term):
    """
    normalize word and remove useless stuff
    """
    return term.lower().\
        replace(":", "").\
        replace(";","").\
        replace(".", "").\
        replace(",","").\
        replace("/","").\
        replace("#","")

def index(filename):
    """
    indexes a given file and saves terms to a posting and non-positional
    inverted index
    """
    global name
    name = filename
    try:
        # open file
        with open(filename, "r") as file:
            docID = 0
            # iterate over each line in file
            for line in file:
                # split them to list of terms
                tweet = line.split()

                for term in tweet:
                    # remove clutter
                    term = normalize(term)
                    # check if term is in stop words to save some memory
                    if term in stop_words:
                        continue
                    # add docID
                    postings[term].add(docID)
                    # update inv_index
                    # we use the term as pointer, because python does not
                    # support pointers, and storing int by indexes will have
                    # massive overhead
                    inv_index[term] = (len(postings[term]), term)
                # increase line number counter
                docID += 1
            # this is for displaying a progress while indexing
    an
```

```

        if docID % 10000 == 0:
            print(str(int(docID / 10000)) + " %")
    except FileNotFoundError as e:
        raise SystemExit("Could not open file: " + str(e))
    return

def getLines(lines):
    result = ""
    try:
        # open file
        with open(name, "r") as file:
            # iterate over the lines and print the lines, which match the
terms
            for i, line in enumerate(file):
                if i in lines:
                    result += str(i) + "\t" + line
    except FileNotFoundError as e:
        raise SystemExit("Could not open file: " + str(e))
    return result

def query(term1, term2=""):
    """
    you can query your search terms. If only one term given it only searches
    for one, otherwise they both have to exist in the tweet
    """
    lines = []
    # remove clutter
    term1 = normalize(term1)
    term2 = normalize(term2)
    # if only one term given look for it
    if term1 in inv_index and not term2:
        (postings_len, postings_pointer) = inv_index[term1]
        # the sorted document_id list out of the postings_list
        lines = sorted(list(postings[postings_pointer]))
    # if two terms are given look for both
    elif term2 and term2 in inv_index:
        (postings_len, postings_pointer) = inv_index[term1]
        # the sorted document_id list out of the postings_list
        lines1 = sorted(list(postings[postings_pointer]))
        (postings_len, postings_pointer) = inv_index[term2]
        # the sorted document_id list out of the postings_list
        lines2 = sorted(list(postings[postings_pointer]))
        # init of iterators
        listiter1 = iter(lines1)
        listiter2 = iter(lines2)
        tmp1 = -1
        tmp2 = -1
        # and intersect the lists to see which lines match both terms
        # if the have the same lines, it will be added to 'lines'
        while True:
            try:
                # like discused in the lesson
                if tmp1 <= tmp2:
                    tmp1 = next(listiter1)
                else:
                    tmp2 = next(listiter2)

```

```
        if tmp1 == tmp2:
            lines.append(tmp1)
        except StopIteration:
            break
    else:
        print("nothing found")
    # we could end here, but we want to get the lines from the file
    return getLines(lines)

if __name__ == '__main__':
    index("tweets")
    print("finished indexing")
    # print(query("geldern"))
    print(query("stuttgart", "bahn"))
```