

Assessment 2

Name: Loukik Bhangale

Reg No: 17BCE0961

Solt: L9+L10

Code:

```
from nltk.corpus import stopwords

from nltk.tokenize import word_tokenize


def cleaner(filename):

    filevariablename=open(filename).read()

    stop_words = set(stopwords.words('english'))

    word_tokens = word_tokenize(filevariablename)

    filtered_sent = [w for w in word_tokens if not w in stop_words]

    filtered_sentence = []

    for w in word_tokens:

        if w not in stop_words:

            filtered_sentence.append(w)

    #print(word_tokens)

    #

    #print("stop words removed!")

    punctuations = list('!"()-[]{};:'"\<>./?@#$$%^&* _~')

    temp=[]

    for char in filtered_sentence:

        if char not in punctuations:

            temp.append(char)

    filtered_sentence=temp

    #print(filtered_sentence)

    return filtered_sentence
```

```

def distinct(doc,li):
    for items in doc:
        if items not in li:
            li.append(items)
doc1=cleaner("Doc 1.txt")
doc2=cleaner("Doc 2.txt")
doc3=cleaner("Doc 3.txt")
doc4=cleaner("Doc 4.txt")
doc5=cleaner("Doc 5.txt")
doc6=cleaner("Doc 6.txt")
doc7=cleaner("Doc 7.txt")
doc8=cleaner("Doc 8.txt")
doc9=cleaner("Doc 9.txt")
doc10=cleaner("Doc 10.txt")
dislist=[]
finaldic={}
posdic={}
def discounter(docnamev,dicvarname):
    temp={}
    temp2={}
    distinct(docnamev,dislist)
    for item in dislist:
        c=docnamev.count(item)
        temp[item]=c
        temp3=[]
        for stuff in range(len(docnamev)):
            if item==docnamev[stuff]:
                temp3.append(stuff)
        temp2[item]=temp3
    finaldic[dicvarname]=temp

```

```

posdic[dicvarname]=temp2
discounter(doc1,"Document 1")
discounter(doc2,"Document 2")
discounter(doc3,"Document 3")
discounter(doc4,"Document 4")
discounter(doc5,"Document 5")
discounter(doc6,"Document 6")
discounter(doc7,"Document 7")
discounter(doc8,"Document 8")
discounter(doc9,"Document 9")
discounter(doc10,"Document 10")
#run it twice because updation distinct list
discounter(doc1,"Document 1")
discounter(doc2,"Document 2")
discounter(doc3,"Document 3")
discounter(doc4,"Document 4")
discounter(doc5,"Document 5")
discounter(doc6,"Document 6")
discounter(doc7,"Document 7")
discounter(doc8,"Document 8")
discounter(doc9,"Document 9")
discounter(doc10,"Document 10")
def printable():
    print("Documents",end="")
    for item in dislist:
        print("\t"+item,end="")
    for item in finaldic:
        print(item,end='\t')
        for stuff in finaldic[item]:
            print(finaldic[item][stuff],end="\t")
    print("\n")

```

#suming up the number of occurance of each words

diclist={}#dictionary of sum of each word occurance in all the documents

for item in dislist:

 t=0

 for stuff in finaldic:

 t+=finaldic[stuff][item]

 diclist[item]=t

#here if you print diclist you get the total occurance of all the words that have appeared in every document used.

def maxword(a):#this finds out if the word exists and if it does then it finds where it occurs maximum and also prints out it's positions

 maxc=0

 maxdoc=0

 if a in dislist:

 for item in finaldic:

 if ((finaldic[item][a])> maxc):

 maxc=finaldic[item][a]

 maxdoc=item

 print ("The word {0} has occured the most in {1}, {2} number of times, in positions:".format(a,maxdoc,maxc))

 print(posdic[maxdoc][a])

 else:

 print ("This word does not appear in the documents")

print ("Enter the number of the function you want to see.\n")

print ("1.See the table of all words and counts respective to their documents.\n")

print ("2.Type in a word and check where it occurs the most and at what positions.\n")

print ("3.See the total occurance of all the words that have appeared in every document used.\n")

n=int(input())

if n ==1:

 printable()

elif n==2:

 t=input()

Sum of occurrence:

jupyter Assessment 2 Last checkpoint: 15 hours ago (unsaved changes) ✓

File Edit View Insert Cell Kernel Widgets Help Snippets Trusted Python 3.0

```
print ("1.See the table of all words and counts respective to their documents.\n")
print ("2.Type in a word and check where it occurs the most and at what positions.\n")
print ("3.See the total occurrence of all the words that have appeared in every document used.\n")
n=int(input())

Enter the number of the function you want to see.

1.See the table of all words and counts respective to their documents.
2.Type in a word and check where it occurs the most and at what positions.
3.See the total occurrence of all the words that have appeared in every document used.

3

In [108]: if n==1:
           printable()
         elif n==2:
           t=input()
           maxword(t)
         elif n==3:
           print(diclist)
         else:
           print("Invalid input.")

{'To': 32, 'Sherlock': 94, 'Holmes': 458, 'always': 59, 'woman': 71, 'I': 3031, 'seldom': 5, 'heard': 113, 'ment
ism': 2, 'name': 62, 'in': 104, 'eyes': 86, 'eclipses': 1, 'predominates': 1, 'whole': 49, 'sea': 1, 'it': 450,
'felt': 32, 'emotion': 5, 'akin': 3, 'love': 18, 'Irene': 16, 'Adler': 15, 'All': 26, 'emotions': 1, 'one': 360,
'particularly': 7, 'abhorrent': 1, 'cold': 21, 'precise': 3, 'admirably': 4, 'balanced': 1, 'mind': 60, 'He': 31
7, 'take': 80, 'perfect': 9, 'reasoning': 11, 'observing': 4, 'machine': 17, 'world': 18, 'seen': 73, 'lower':
7, 'would': 335, 'placed': 19, 'false': 4, 'position': 23, 'never': 85, 'spoke': 26, 'softer': 1, 'passions': 1,
'save': 36, 'gibe': 1, 'sneer': 5, 'They': 63, 'admirable': 3, 'things': 33, 'observer': 4, '---': 231, 'excellen
t': 15, 'drawing': 6, 'well': 8, 'men': 48, 's': 372, 'motives': 4, 'actions': 4, 'But': 178, 'trained': 4, 're
asoner': 6, 'admit': 3, 'intrusions': 1, 'delicate': 7, 'finely': 1, 'adjusted': 2, 'temperament': 1, 'introduce
': 5, 'distracting': 1, 'factor': 3, 'might': 124, 'throw': 11, 'doubt': 65, 'upon': 479, 'mental': 1, 'results
': 14, 'Gut': 1, 'sensitive': 1, 'instrument': 3, 'crack': 3, 'high-power': 1, 'lenses': 1, 'disturbing': 1, 'a
strong': 33, 'nature': 23, 'had': 191, 'yet': 78, 'late': 30, 'dubious': 1, 'questionable': 2, 'memory': 12, 'lit
tle': 274, 'lately': 7, 'My': 92, 'marriage': 26, 'drifted': 2, 'us': 192, 'away': 109, 'complete': 13, 'happine
ss': 3, 'home-centred': 1, 'interests': 4, 'rise': 6, 'around': 3, 'man': 305, 'first': 100, 'finds': 4, 'master
': 13, 'establishment': 1, 'sufficient': 7, 'absorb': 1, 'attention': 25, 'loathed': 1, 'every': 84, 'form': 21,
'society': 15, 'Bohemian': 4, 'soul': 11, 'remained': 20, 'lodgings': 8, 'Baker': 42, 'Street': 61, 'buried': 6,
'among': 43, 'old': 70, 'books': 11, 'alternating': 1, 'weak': 27, 'cocaine': 3, 'ambition': 1, 'drowsiness': 1,
'drug': 4, 'fierce': 6, 'energy': 7, 'seen': 14, 'still': 70, 'ever': 72, 'deeply': 11, 'attracted': 5, 'study':
13, 'crime': 31, 'occupied': 1, 'immense': 10, 'faculties': 2, 'extraordinary': 20, 'powers': 8, 'observation':
```