impost nitk MEX. down ( Woodnet) NIEL. download ( Puncti) text = "This is Andrew's fext, isn't Pt?" 1) - How many . to kens: tokenizer = M+K. tokenize. Who tespacetokenizer () tokens = tokenizex. tokenize (text) print (len (tobens)) print (tokens) 2) How many tokens; transancherd tokenious? tokenfrer = nitk. fokenire. Weiteracontokonizar () tokens: tokeniar. tokenizer ('text) print ( lon ( tokens)) print (bolans) 3) How many tokens there are & you we Woodfurd. Tokenizex? tokanizer = MIK, tokanize, Wordfundt Tokarizer() forens = foren (zer-forentze (text) evint (cen (tokens)) tight (tokens)

# Natural Language Processing Lab Lab1. Understanding Large Text Files

### **EXERCISE-1**

-

9

-

-

Consider the following text.

```
import nltk
nltk.download('wordnet')
text = "This is Andrew's text, isn't it?"
```

 $1. \ \ \text{How many tokens are there if you use WhitespaceTokenizer?}. \ \ \text{Print tokens.}$ 

```
tokenizer = nltk.tokenize.WhitespaceTokenizer()
tokens = tokenizer.tokenize(text)
print(len(tokens))
print(tokens)
```

2. How many tokens are there if you use TreebankWordTokenizer?. Print tokens.

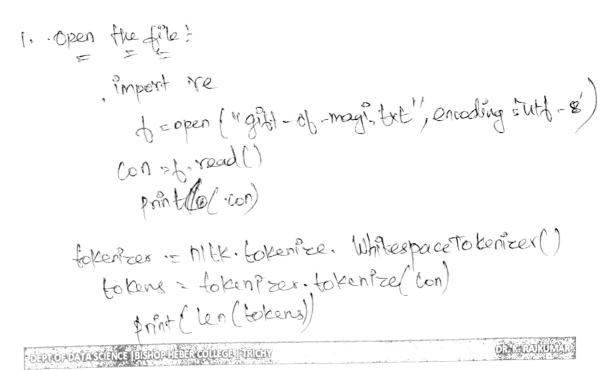
```
tokenizer = nltk.tokenize.TreebankWordTokenizer()
```

3. How many tokens there are if you use WordPunctTokenizer?. Print tokens.

tokenizer = nltk.tokenize.WordPunctTokenizer()

#### **EXERCISE-2**

- 1. Open the file: O. Henry's The Gift of the Magi (gift-of-magi.txt).
- 2. Write a Python script to print out the following:
  - 1. How many word tokens there are
  - 2. How many word types there are, (word types are a unique set of words)
  - 3. Top 20 most frequent words and their counts
  - 4. Words that are at least 10 characters long and their counts
  - 5. 10+ characters-long words that occur at least twice, sorted from most frequent to least



```
dram nitt Propost
  data = Freg Dist (tokens)
  dala.
  data. most_common (20)
A) from NIX Emport.
   text = [w forw in tokens is len(w)>10]
    Print (text)
    freq = trappist(text)
     freq
     boext = [w for w in tokens is len(w) >10]
        3 = Frag Dist (text)
    der i, i is freq. Ptem():
         3/ lan (P)) 10 and 3>21
           print(2,3)
  ExexCise - 3
       Fist Comprevension
  step-1
       frame: 11. Idata austen-emma. txt"
                                                        C
         d = 6pen ("austen-emma-txt", encoding: 14-8)
        etxt= firead()
        f. Close()
      etxt[-200;]
```

# **EXERCISE -3: List Comprehension**

# STEP-1

Download the document Austen's *Emma* ("austen-emma.txt"). Read it in and apply the usual text processing steps, building three objects: etoks (a list of word tokens, all in lowercase), etypes (an alphabetically sorted word type list), and efreq (word frequency distribution).

```
>>> fname = "./data/austen-emma.txt"
>>> f = open(fname, 'r')
>>> etxt = f.read()
>>> f.close()
>>> etxt[-200:]
'e deficiencies, the wishes,\nthe hopes, the confidence, the predictions of the
small band\nof true friends who witnessed the ceremony, were fully answered\nin
the perfect happiness of the union.\n\nFINIS\n'
>>> etoks = nltk.word_tokenize(etxt.lower())
>>> etoks[-20:]
['of', 'true', 'friends', 'who', 'witnessed', 'the', 'ceremony', ',', 'were',
'fully', 'answered', 'in', 'the', 'perfect', 'happiness', 'of', 'the', 'union',
'.', 'finis']
>>> len(etoks)
191781
>>> etypes = sorted(set(etoks))
>>> etypes[-10:]
['younger', 'youngest', 'your', 'yours', 'yourself', 'yourself.', 'youth', 'youthful',
'zeal', 'zigzags']
>>> len(etypes)
7944
>>> efreq = nltk.FreqDist(etoks)
>>> efreq['beautiful']
24
```

## STEP 2: list-comprehend Emma

Now, explore the three objects wlist, efreq, and etypes to answer the following questions. Do NOT use the for loop! Every solution must involve use of LIST COMPREHENSION.

# Question 1: Words with prefix and suffix

What are the words that start with 'un' and end in 'able'?

### Question 2: Length

-

How many Emma word types are 15 characters or longer? Exclude hyphenated words.



tokenizer = nth. tokenize. Whitesparetokenizer() tokens = , Lokenfoor Lokenfoe (txt) totens [-20] etacs = nitc. word\_tokenite(etxt.lower()) etoks[-20:] len (etaks) etypes: sorted (set(ebbs)) etypes [-10:] lan (etypes) educy = nitk. Freq DBf(etoks) etra ! beautifu! Step-2: Question 1: Words with prefix and suffix: [word for food in tokens . it word . etaxts with ["un") \$ word, end suith ("able") Question 2 = tergth How many. Emma word types are 15 characters of longer. to Kentrer := ntk. fokenitee. libral Runt tokenizer() for stokenizer to kenize letat) [ Word for word in take if lan (word) 215 ?

Question 3: Average word length

What's the average length of all Emma word types?

Question 4: Word frequency

How many Emma word types have a frequency count of 200 or more? How many word types appear only once?

Question 5: Emma words not in wlist

Of the Emma word types, how many of them are not found in our list of ENABLE English words, i.e., wlist?

## STEP 3: bigrams in Emma

Let's now try out bigrams. Build two objects: e2grams (a list of word bigrams; make sure to cast it as a list) and e2gramfd (a frequency distribution of bigrams) as shown below, and then answer the following questions.

>>> e2grams = list(nltk.bigrams(etoks))

>>> e2gramfd = nltk.FreqDist(e2grams)

>>>

Question 6: Bigrams

What are the last 10 bigrams?

Question 7: Bigram top frequency

What are the top 20 most frequent bigrams?

Question 8: Bigram frequency count

How many times does the bigram 'so happy' appear?

Question 9: Word following 'so'

What are the words that follow 'so'? What are their frequency counts? (For loop will be easier; see if you can utilize list comprehension for this.)

**Question 10: Trigrams** 

What are the last 10 trigrams? (You can use nltk.util.ngrams() method)

Question 11: Trigram top frequency

What are the top 10 most frequent trigrams?

Question 12: Trigram frequency count

How many times does the trigram 'so happy to' appear?

What's the average length of all Emma Word types?

average = Sum (len (word) for word in toke) / len (toke)

DEPT OF DATA SCIENCE BISHOP HEBER COLLEGE LITRICHY

lg = [J for i' in toke! Po len(2) 315! 19. append (?) print (la) Greeny 4: Word frequency How many Emma Word types have a fragmency bount. . from Nith . Import. falemn = Freq Dest (toke) for is. in fdiemm. (tem(): Exist (6,3) How many lord types appears only once? for P, g. in foremm. Plans(), proint (P, 9). What are the last 10 bigrams. ez granis = 12st (netk. bigranus (toke)) eagrams |d = Mtk. Freq Dest (02 growns) ezgramfd. 0 last-ten: Freq Dest (dect (eggram folamost comment) [-10:] last-ten

**NOTES** Question 7: Blysam top frequency. to Kenizer= nltk. to konize. White space tokenieur () toker = fokenizer. fokenize (etxt) ezgrams: l'est (nith. bigrams (fokens)) exgrampd = nith. Frag Dist (exgrams) ezgranfol. most-comma (20) Question 8 - Bigram frequency count. How many firmer does the Digrame ! so happy appear? for i, g in ezgeramd. items(): for (= = [ ( so / happy )): beint ( ))) Question 9: Word following 1 so from collections-impost-counter. Import re Words = re-findall ( r'sot / W+, open l'ausen-enma. Ext), read()) ab = counter, (zip (words)) point (ab) Question 10. e. agrans = list (nltk. bylgsams (tokens)) esgranded = NHK. Freq Dist (e3grams) DEPLOE DATASCIENCE BISHOP HEBER COLLEGE | MRICHY

( S)

( CO.

63

30

200

-

-840

3

-3

3

3