

NLP Homework 2

Sentiment Analysis: Exploratory Analysis

1. Introduction

In this homework we will perform exploratory and sentiment analysis on the same corpora of Shakespearean books used in HW-1.

Below is an excerpt showing indexing books from the Gutenberg corpora:

```
1 # HW-2
2 # Exploratory exercise for sentiment analysis
3 # Finding adverb and adjective phrases, and computing basic statistics
4
5 # importing required nltk libraries
6 import nltk
7 from nltk import sent_tokenize
8
9 # In HW-1, I used books written by Shakespeare – Caesar and Hamlet. We will continue exploratory
10 # sentiment analysis on the same books
11 nltk.corpus.gutenberg.fileids()
12
13 # Get Shakespeare books in the Gutenberg corpus
14 shakespeare_books = [book for book in nltk.corpus.gutenberg.fileids() \
15                       if 'shakespeare' in book]
16
17 # Book-1: Caesar (Genre: Tragedy)
18 caesar = nltk.corpus.gutenberg.raw(shakespeare_books[0])
19
20 # Book-2: Hamlet (Genre: Tragedy/Comedy)
21 hamlet = nltk.corpus.gutenberg.raw(shakespeare_books[1])
22
23 print(caesar[:50])
24 print(hamlet[:50])
```

```
[The Tragedie of Julius Caesar by William Shakespe
[The Tragedie of Hamlet by William Shakespeare 159
```

2. Processing the corpus

The HW will adopt the existing code from *HMW 2-Code Ideas.ipynb* to the above corpora and provide a detailed code walk through and explanation of all steps used to answer the following questions:

- List top-50 adjective/adverb phrases (by frequency)
- List top-50 adjective/adverb words
- List top-50 nouns or verbs
- Average length of sentence

First, we will split the two corpora into sentences using `nltk.sent_tokenize()` before iterating and word tokenizing sentences using `nltk.word_tokenize()`

```
1 # Apply the word tokenizer to each sentence
2 token_caesar = [nltk.word_tokenize(sent) for sent in caesar_split]
3 print(token_caesar[:2])
4
5 # the output is a list of strings that contains the sentences
6 print('Type %s' % (type(token_caesar)))
7 print('Caesar tokens: %d' % (len(token_caesar)))
8
9 # Repeat the same for hamlet text
10 token_hamlet = [nltk.word_tokenize(sent) for sent in hamlet_split]
11 print(token_hamlet[:2])
12 print('Hamlet tokens: %d' % (len(token_hamlet)))
```

Next, we will use the Stanford POS tagger to tag tokens:

```
1 ## POS Tagging, to retrieve adjective (JJs) and adverb (RBs) tags
2
3 # use the Stanford POS tagger to POS tag tokens of each sentence
4 # this is the default tagger in nltk
5 caesar_tagged = [nltk.pos_tag(tokens) for tokens in token_caesar]
6 print(caesar_tagged[:2])
7
8 hamlet_tagged = [nltk.pos_tag(tokens) for tokens in token_hamlet]
9 print(hamlet_tagged[:2])
```

Output:

```
[(['', 'IN'), ('The', 'DT'), ('Tragedie', 'NNP'), ('of', 'IN'), ('Julius', 'NNP'), ('Caesar', 'NNP'), ('by', 'IN'), ('William', 'NNP'), ('Shakespeare', 'NNP'), ('1599', 'CD'), (']', 'NNP'), ('Actus', 'NNP'), ('Primus', 'NNP'), ('.', '.')], [('Scoena', 'NNP'), ('Prima', 'NNP'), ('.', '.')]]
[(['', 'IN'), ('The', 'DT'), ('Tragedie', 'NNP'), ('of', 'IN'), ('Hamlet', 'NNP'), ('by', 'IN'), ('William', 'NNP'), ('Shakespeare', 'NNP'), ('1599', 'CD'), (']', 'NNP'), ('Actus', 'NNP'), ('Primus', 'NNP'), ('.', '.')], [('Scoena', 'NNP'), ('Prima', 'NNP'), ('.', '.')]]
```

2.1 List the top 50 adjective/adverb phrases (by frequency)

For this purpose, we will use the technique of chunking. NLTK offers two utilities, *RegexpParser()*, to define or update custom grammar rules and *parse()* to chunk and parse through the sentence.

Since the above was going to be used across corpora and different POS tags I decided to write it as a python function as seen below:

```
10 def chunking(grammar, taggedtext, metadata):
11     # Extract the metadata
12     label = metadata['label']
13     desc = metadata['desc']
14     text = metadata['text']
15
16     # Second step: import the nltk parser to process each sentence
17     chunk_parser = nltk.RegexpParser(grammar)
18
19     _tags = []
20     for sent in taggedtext:
21         if len(sent) > 0:
22             tree = chunk_parser.parse(sent)
23             for subtree in tree.subtrees():
24                 if subtree.label() == label:
25                     _tags.append(subtree)
26
27     # Visualizing the actual adj/adv phrase
28     _phrases = []
29     for sent in _tags:
30         temp = ''
31         for w, t in sent:
32             temp += w + ' '
33         _phrases.append(temp)
34
35     print('First 10 %s phrases (%s): %s' % (desc, text, _phrases[:10]))
36
37     # Following our NLTK textbook, chapter 1 on Language Processing (https://www.nltk.org/book/ch01.html)
38     # FREQUENCY DISTRIBUTIONS
39     # Top 50 adjective phrases
40     _freq = nltk.FreqDist(_phrases)
41
42     print('Top %s phrases by frequency (%s): ' % (desc, text))
43     for word, freq in _freq.most_common(50):
44         print(word, freq)
45
46     #print the list of our sentences:
47     print('Length of %s phrase sentences (%s): %d' % (desc, text, len(_tags)))
48
49     return _tags
50
```

In this function, we would:

- Define a parser based on custom grammar which is one of the args
- Iterate through the sentences in the corpora and use the parser
- Walk the subtrees and match for labels like ADJPH, ADVPH etc
- And last, walk `_tags` list and collect phrases to run through the `most_common()` function to mine for frequencies

Below is an excerpt of code showing an instance of calling or invoking the `chunking()` function for the adjective phrase:

```
52 grammar_adjph = "ADJPH: {<RB.?>+<JJ.?>}"
53 # This regex reads as: "find groups ("< >") of RBs (adverbs) together with groups of JJs (adjectives), with grou
54 # RBs with any ending (the "." is a placeholder or wildcard for the "R" and the "S" at the end of RBR and RBS,
55 # while "?" indicates "optional character" so RB can be found alone as well). Same regex operators apply to JJs.
56
57 caesar_adjph_tags = chunking(grammar_adjph, caesar_tagged, {'text': 'Caesar', 'label': 'ADJPH', 'desc': 'adjective'
58 print('-----')
59 hamlet_adjph_tags = chunking(grammar_adjph, hamlet_tagged, {'text': 'Hamlet', 'label': 'ADJPH', 'desc': 'adjective'
```

Output:

First 10 adjective phrases (Caesar): ['Truly sir ', 'Truly sir ', 'then senslesse ', 'there haue ', 'most exalted ', 'not mou ', 'yet againe ', 'once againe ', 'too stubbor ne ', 'too strange ']

Top adjective phrases by frequency (Caesar):

```
so much 7
too much 4
so great 3
then thy 3
so good 3
Truly sir 2
so many 2
most Noble 2
more worthy 2
not backe 2
not meete 2
```

Length of adjective phrase sentences (Caesar): 143

First 10 adjective phrases (Hamlet): ['now strook ', 'once againe ', 'So frown ', 'most observant ', 'most emulate ', 'So hallow ', 'so gracious ', 'So haue ', 'so farre ', 'not fayl ']

Top adjective phrases by frequency (Hamlet):

```
so farre 5
too much 3
as much 3
too blame 3
So much 2
not thy 2
not fit 2
so much 2
most excellent 2
```

Length of adjective phrase sentences (Hamlet): 237

Next, we will call the *chunking()* function for Adverb phrases:

```
1 # Now we look for "adverb phrases" or chunks that have 2 consecutive adverbs ('RB')
2 # First step: writing a grammar that defines POS rules of the adverb phrase the chunk
3 # we name this grammar "ADVPH" ("ADVerb PHrase")
4 grammar_advph = "<RB>+<RB>+"
5
6 caesar_advph_tags = chunking(grammar_advph, caesar_tagged, {'text': 'Caesar', 'label': 'ADVPH', 'desc': 'adverb'})
7 print('-----')
8 hamlet_advph_tags = chunking(grammar_advph, hamlet_tagged, {'text': 'Hamlet', 'label': 'ADVPH', 'desc': 'adverb'})
```

Output:

First 10 adverb phrases (Caesar): ['art not ', 'So well ', 'heere so long ', 'as well ', 'as well ', 'so indeed ', 'till now ', 'not so ', 'So soone ', 'not then ']

Top adverb phrases by frequency (Caesar):

```
not so 9
not well 4
as well 3
heere so 3
So well 2
art not 1
heere so long 1
so indeed 1
till now 1
So soone 1
not then 1
```

Length of adverb phrase sentences (Caesar): 74

First 10 adverb phrases (Hamlet): ['spoke too ', 'Thus twice before ', 'doth well ', 'Thus much ', 'Not so ', 'too too ', 'not so much ', 'too roughly ', 'not well ', 'else neere ']

Top adverb phrases by frequency (Hamlet):

```
not so 5
not well 4
thee well 2
very well 2
So much 2
too much 2
so well 2
spoke too 1
Thus twice before 1
doth well 1
```

Length of adverb phrase sentences (Hamlet): 84

2.2 List the top-50 adjective/adverb words

For this exercise, we would define a function *top_tokens()* taking POS tags as one of the arguments (type list).

```

1 # Top 50 tokens (grouped by Adjective, Adverb or Nouns)
2
3 def top_tokens(taggedtext, pos_list):
4     _tokens = []
5     for sentence in taggedtext:
6         for word, pos in sentence:
7             if pos in pos_list:
8                 if len(word)>1:
9                     _tokens.append(word)
10    freq_pos = nltk.FreqDist(_tokens)
11
12    for word, freq in freq_pos.most_common(50):
13        print(word, freq)
14
15 # Top 50 adjective tokens
16 print('Top 50 Adjective tokens:')
17 top_tokens(caesar_tagged, ['JJ', 'JJR', 'JJS'])
18 top_tokens(hamlet_tagged, ['JJ', 'JJR', 'JJS'])

```

Within the function, we would:

- Iterate the POS tag tokens and find a match in *pos_list*
- If there's a match, we append it to a list *_tokens* and then run an *nltk.FreqDist()* utility and use the *most_common()* words to mine for frequency

In the excerpt above, we see the usage for Adjective tokens and the corresponding output below.

Output:

```

Adjectives (Caesar):
good 48
thy 41
Noble 32
great 25
thou 24
such 23
much 22
true 18
Good 18
many 15
dead 14

```

```

Adjectives (Hamlet):
good 76
thy 54
more 37
such 34
most 30
much 25
dead 25
true 21
Good 21
thou 20

```

Next, using the same *top_tokens()* call, we mine for top-50 adverbs:

```

1 # Top 50 adverb tokens
2 print('Top 50 Adverb tokens:')
3 top_tokens(caesar_tagged, ['RB', 'RBR', 'RBS'])
4 top_tokens(hamlet_tagged, ['RB', 'RBR', 'RBS'])

```

Output:

Adverbs (Caesar):

```
not 255
so 103
then 79
well 40
now 39
too 30
Then 28
yet 27
heere 26
So 24
more 24
Now 24
there 15
once 13
thus 13
```

Adverbs (Hamlet):

```
not 313
so 139
then 75
now 68
well 53
too 50
more 46
very 44
most 35
So 33
Then 33
thus 33
Now 24
yet 23
```

And finally, nouns and verbs using the same custom function:

```
1  ## TO DO / YOUR TURN NOW!
2  ## NOUN EXTRACTION
3  ## VERB EXTRACTION
4  ## REMEMBER TO CHECK THE PENN POS TAGS LIST: https://www.ling.upenn.edu/courses/Fall\_2003/ling001/penn\_treebank\_
5  ## TO FIND ALL TAGS
6
7  print('Top 50 Noun tokens:')
8  top_tokens(caesar_tagged, ['NN', 'NNS', 'NNP', 'NNPS']) #Noun, Noun-plural, Noun-Propor, Noun-Propor-plural
9  top_tokens(hamlet_tagged, ['NN', 'NNS', 'NNP', 'NNPS'])
10
11
12  print('\nTop 50 Verb tokens:')
13  top_tokens(caesar_tagged, ['VB', 'VBD', 'VBG', 'VBP', 'VBZ']) # Verb, Verb-past-tense, Verb-present participle,
14  # Verb-past participle, singular present (non-3rd)
15  # singular present (3rd)
16  top_tokens(hamlet_tagged, ['VB', 'VBD', 'VBG', 'VBP', 'VBZ'])
```

Output:

Nouns (Caesar):

```
Caesar 187
Brutus 160
Bru 152
```

```
Cassi 107
Cassius 85
Antony 75
Enter 58
men 57
man 54
thou 49
Ant 48
Lord 44
```

```
Nouns (Hamlet):
Ham 334
Lord 210
King 170
Hamlet 98
Hor 95
Enter 80
Qu 62
Laer 59
Ile 58
Ophe 55
Pol 48
```

```
Verbs (Caesar):
is 247
be 132
do 107
haue 100
are 96
was 64
know 63
did 61
am 52
let 41
```

```
Verbs (Hamlet):
is 349
be 175
haue 129
are 111
do 79
was 79
know 66
's 65
let 59
come 50
```

2.3 Average length of sentences

The average length of sentences can be computed by summing the length of all sentences and dividing it by the number of sentences in *caesar_split* or *hamlet_split* (which is essentially output of `nltk.sent_tokenize` function).

Below is a code excerpt:


```

1 # Following our NLTK textbook, Writing Structural Programs chapter
2 # section on Procedural vs Declarative style (http://www.nltk.org/book\_1ed/ch04.html)
3
4 ## CORPUS STATISTICS--SENTENCES LENGTH
5
6 # Calculating the average length of sentences in the entire corpus
7 # from http://www.nltk.org/book\_1ed/ch04.html
8 caesar_total_corpus = sum(len(sent) for sent in caesar_split)
9 print('Average len of sentence (Caesar): %s' %(caesar_total_corpus / len(caesar_split)))
10
11 hamlet_total_corpus = sum(len(sent) for sent in hamlet_split)
12 print('Average len of sentence (Caesar): %s' %(hamlet_total_corpus / len(hamlet_split)))

```

Output:

```

Average len of sentence (Caesar): 69.16394472361809
Average len of sentence (Hamlet): 67.78853503184713

```

3. Conclusion

As we make further in-roads into the corporal statistics using NLP let us review data analyzed in this homework for both the Caesar and Hamlet Shakespearean texts.

Here's a summary:

Statistics	Caesar corpus	Hamlet corpus	Comment
Length of adjective phrases in corpus	143	237	As much as 1.65x more adjective phrases used in the Hamlet corpus
Length of adverb phrases in corpus	74	84	Marginal difference in the number of adverb phrases
Average length of sentence	69	67	Almost same length of sentences

As seen through each of the sections in this homework for adjective phrases, adverb phrases, adjectives, adverbs, nouns, or verbs very subtle differences are seen. The language or style adopted for each of these Shakespearean plays for most part, seem comparable.