

Spam Content - Report

Group 26:

- Samarjoy Pandit
- Saksham
- Vishal
- Akanksha Jojwan
- Saurabh

Source Code:

1. Count articles (a, an, the)

```
1 text = open('sample.txt').read()
2 a_count = text.count('a') + text.count('A')
3 an_count = text.count('an') + text.count('An')
4 the_count = text.count('the') + text.count('The')
5 print('Number of a:', a_count)
6 print('Number of an:', an_count)
7 print('Number of the:', the_count);
```

Screenshot:

Average length of words in the file: 5.413990982374641

2. Average length of words

```
1 s = open('sample.txt').read()
2 numlist = list(map(len, s.split()))
3 sum = 0
4 for x in numlist:
5     sum = sum + x
6 print('Average length of words in the file:', sum/len(numlist))
```

Screenshot:

Number of a: 3178
Number of an: 566
Number of the: 469

3. tri-grams

```

1 import nltk
2 def get_words(string):
3     tokenizer = nltk.tokenize.RegexpTokenizer(r'\w+')
4     return tokenizer.tokenize(string)
5 string = open('sample.txt').read()
6 words = get_words(string)
7 finder = nltk.collocations.TrigramCollocationFinder.from_words(words)
8 scored = finder.score_ngrams(nltk.collocations.TrigramAssocMeasures().raw_freq)
9 trigrams = finder.ngram_fd.items()
10 for x in trigrams:
11     print(x+"\n")

```

Screenshots:

```

(('Banner', 'logo', 'This'), 1)
(('logo', 'This', 'November'), 1)
(('This', 'November', 'is'), 1)
(('November', 'is', 'the'), 1)
(('is', 'the', 'Wikipedia'), 1)
(('the', 'Wikipedia', 'Asian'), 1)
(('Wikipedia', 'Asian', 'Month'), 1)
(('Asian', 'Month', 'Come'), 1)
(('Month', 'Come', 'join'), 1)
(('Come', 'join', 'us'), 1)

```

```

(('Search', 'engine', 'optimization'), 6)
(('engine', 'optimization', 'From'), 2)
(('optimization', 'From', 'Wikipedia'), 2)
(('From', 'Wikipedia', 'the'), 2)
(('Wikipedia', 'the', 'free'), 2)
(('the', 'free', 'encyclopedia'), 2)
(('free', 'encyclopedia', 'SEO'), 2)
(('encyclopedia', 'SEO', 'redirects'), 2)
(('SEO', 'redirects', 'here'), 2)

```