

```
1 # step 1 : Loading and reading the Data using os.path
2 import os
3 import sys
4 with open(os.path.join(sys.path[0], "ID1(reading and
  its correlates2014).txt"), "r") as f:
5     f.close()
6
7
8 # step 2 :Split by Whitespace and preserve the
  punctuations.
9 file=open("ID1(reading and its correlates2014).txt"
  , "r")
10 text=file.read()
11 words=text.split()
12 print(words[:])
13 #punctuation is preserved (e.g. "wasn't" and "armour-
  like"),
14 # which is nice. We can also see that end of
15 #sentence punctuation is kept with the
16 # last word (e.g. "thought.")
17
18
19 # Step3 :Split by Whitespace and obliterate the
  punctuations.
20 #want the words, but without the punctuation like
  commas
21 #and quotes. We also want to keep contractions
  together.
22 #load the text file, split it into words by white
  space,
23 #then translate each word to remove the punctuation.
24 import string
25 table= str.maketrans("", "", string.punctuation)
26 stripped = [w.translate(table) for w in words]
27 print(stripped[:])
28
29 #printing huge number of punctuations in python
30 print(string.punctuation)
31
32
33
34 #Normalizing Case
35 words=[word.lower() for word in words]
36 print(words[:])
```

```
37
38 words=[word.capitalize() for word in words]
39 print(words[:])
40
41
42
43 print("An end to this task; lets go for NLTK")
44
45
46
47 # Remember, the simpler, the better.
48 # Simpler text data, simpler models, smaller
49 # vocabularies. You can always make things
50 # more complex later to see if it results in
51 # better model skill
52 # we'll look at some of the tools in the NLTK (a
   Python library
53 # written for working and modeling text.)
54 # library that offer more than simple string
   splitting.
55 # NLTK = Natural Language Tool-Kit
56 import nltk
57 #nltk.download()
58 filename = 'ID1(reading and its correlates2014).txt'
59 file= open(filename, "r")
60 text= file.read()
61 print(text)
62 file.close()
63 #NLTK provides the sent_tokenize() function to split
64 #text into sentences.
65 from nltk.tokenize import sent_tokenize
66 sentences = sent_tokenize(text)
67 print(sentences[5])
68 total_sentences=len(sentences)
69 print(total_sentences)
70 for sentence in range (total_sentences):
71     print("Sentence ",sentence, "is",": ",sentences[
        sentence])
72 #NLTK provides a function called word_tokenize()
73 #for splitting strings into tokens (nominally words).
74 from nltk.tokenize import word_tokenize
75 words = word_tokenize(text)
76 print(words[51])
77 total_words=len(words)
```

```

78 print(total_words)
79 for word in range(total_words):
80     print("word ",word,"is",words[word])
81
82 #Our next step is to filtering out the punctuation.
83 filename="ID1(reading and its correlates2014).txt"
84 file=open("ID1(reading and its correlates2014).txt"
85           ,"r")
86 text= file.read()
87 from nltk.tokenize import word_tokenize
88 words = word_tokenize(text)
89 only_words = [word for word in words if word.isalpha
90               ()]
91 print(only_words[:])
92
93 #how to remove stopwords?
94 print("the task tries to eliminate the words
95       extracted "
96       "here called stop words")
97 from nltk.corpus import stopwords
98 stop_words=stopwords.words("english")
99 print(stop_words)
100 #lets compare your tokens to the stop words and
101     filter them out
102 print("lets creat a corpus with no stop words")
103 words=[word for word in words if not word in
104        stop_words]
105 print(words[:])
106
107 #lets analyze for stem of the corpus with no stop
108     words.
109 file_name = "ID1(reading and its correlates2014).
110             txt"
111 file= open("ID1(reading and its correlates2014).txt
112            ","rt")
113 text=file.read()
114 from nltk.tokenize import word_tokenize
115 tokens= word_tokenize(text)
116 total_token=len(tokens)
117 print(total_token)
118 for words in range (total_token):
119     print("The word number ", words, " is ", tokens[

```

```
113 words])
114
115 from nltk.stem.porter import PorterStemmer
116 porter= PorterStemmer()
117 stemmed=[porter.stem(word) for word in tokens]
118 print(stemmed)
119 total_stem=len(stemmed)
120 print(total_stem)
121 for item in range (total_stem) :
122     print("The stem, number ", item, " is ",stemmed[
        item])
123
124
125 # Extracting Method (lets you take a code fragment
126 # that can be grouped, move it into a separated
127 # method, and replace the old code with a call to
128 the method.
129 # When you extract the method you need to check for
130 # variables. If there is one output variable, it
131 # is used as a return value for the extracted
132 # method. In case there are multiple output
133 # variables, the Extract Method refactoring may
134 # not be applied, and the error message appears.
135
136
137
```