

Document similarity will consist of three fundamental steps:

Split the documents in words. Compute the word frequencies. Calculate the dot product of the document vectors. **bold text**

```
1 import math
2 import string
3 import sys
4
5 # reading the text file
6 # This function will return a
7 # list of the lines of text
8 # in the file.
9 def read_file(filename):
10
11     try:
12         with open(filename, 'r') as f:
13             data = f.read()
14             return data
15
16     except IOError:
17         print("Error opening or reading input file: ", filename)
18         sys.exit()
19
20 # splitting the text lines into words
21 # translation table is a global variable
22 # mapping upper case to lower case and
23 # punctuation to spaces
24 translation_table = str.maketrans(string.punctuation+string.ascii_uppercase,
25                                   " "*len(string.punctuation)+string.ascii_lowercase)
26
27 # returns a list of the words
28 # in the file
29 def get_words_from_line_list(text):
30
31     text = text.translate(translation_table)
32     word_list = text.split()
33
34     return word_list
35
```

Now that we have the word list, we will now calculate the frequency of occurrences of the words.

```
1 # counts frequency of each word
2 # returns a dictionary which maps
3 # the words to their frequency.
4 def count_frequency(word_list):
5
6     D = {}
7
```

```

8   for new_word in word_list:
9
10      if new_word in D:
11          D[new_word] = D[new_word] + 1
12
13      else:
14          D[new_word] = 1
15
16   return D
17
18 # returns dictionary of (word, frequency)
19 # pairs from the previous dictionary.
20 def word_frequencies_for_file(filename):
21
22   line_list = read_file(filename)
23   word_list = get_words_from_line_list(line_list)
24   freq_mapping = count_frequency(word_list)
25
26   print("File", filename, ":", )
27   print(len(line_list), "lines, ", )
28   print(len(word_list), "words, ", )
29   print(len(freq_mapping), "distinct words")
30
31   return freq_mapping
32

```

Lastly, we will calculate the dot product to give the document distance.

```

1 # returns the dot product of two documents
2 def dotProduct(D1, D2):
3     Sum = 0.0
4
5     for key in D1:
6
7         if key in D2:
8             Sum += (D1[key] * D2[key])
9
10    return Sum
11
12 # returns the angle in radians
13 # between document vectors
14 def vector_angle(D1, D2):
15     numerator = dotProduct(D1, D2)
16     denominator = math.sqrt(dotProduct(D1, D1)*dotProduct(D2, D2))
17
18     return math.acos(numerator / denominator)
19

```

```

1 def documentSimilarity(filename_1, filename_2):
2
3
4     sorted_word_list_1 = word_frequencies_for_file(filename_1)
5     sorted_word_list_2 = word_frequencies_for_file(filename_2)

```

```
5 sorted_word_list_2 = word_frequencies_for_file(filename_2)
6 distance = vector_angle(sorted_word_list_1, sorted_word_list_2)
7
8 print("The distance between the documents is: % 0.6f (radians)"% distance)
9
```

```
1 # Driver code
2 documentSimilarity('alice.txt', 'alicemodified.txt')
```

```
📄 File alice.txt :
1326 lines,
257 words,
138 distinct words
File alicemodified.txt :
1335 lines,
259 words,
139 distinct words
The distance between the documents is: 0.042916 (radians)
```