

## IT - 425 Natural Language Processing

Q1) a)

	Python	Theory	Geometry	Worked	Research
doc1	1	0	0	0	0
doc2	0	1	1	1	1
doc3	0	1	0	1	1
doc4	0	0	0	1	0
doc5	1	0	0	0	0
doc6	0	0	0	0	1

We consider Research as our text (test) vector.

Now,

$$\text{Cosine similarity } (u, v) = \frac{u \cdot v}{|u| |v|}$$

$$\text{where } |u| = \sqrt{\sum_{i=1}^n u_i^2}$$

$$\begin{aligned} \text{Cosine (python, Research)} &= \frac{(1000 \ 10)^T (0 \ 1 \ 0 \ 0 \ 1)}{|1000 \ 10| |0 \ 1 \ 0 \ 0 \ 1|} \\ &= 0 \end{aligned}$$

$$\text{cosine}(\text{theory, research}) = \frac{(0, 1, 1, 0, 0, 0)^T (0, 1, 1, 0, 0, 1)^T}{(0, 1, 1, 0, 0, 0) (0, 1, 1, 0, 0, 1)}$$

$$\frac{1+1}{\sqrt{\sum u_i^2} \sqrt{\sum v_i^2}} = \frac{2}{(|u|/|v|)}$$

$$|u| = \sqrt{1+1} = 1.414$$

$$|v| = \sqrt{3} = 1.73$$

$$\text{CS}(\text{theory, research}) = \frac{2}{(1.41)(1.73)} = 0.8164$$

$$\text{CS}(\text{geometry, research}) = \frac{(0, 1, 0, 0, 0)^T (0, 1, 0, 0, 1)^T}{(0, 1, 0, 0, 0) (0, 1, 0, 0, 1)}$$

$$= \frac{1}{\sqrt{1} \sqrt{3}} = \frac{1}{1.73} = 0.5773$$

$$\text{cosine sim}(\text{worked, research}) = \frac{(0, 1, 1, 0, 0)^T (0, 1, 1, 0, 0)^T}{(0, 1, 1, 0, 0) (0, 1, 1, 0, 0)}$$

$$= \frac{2}{\sqrt{3} \sqrt{3}} = \frac{2}{3} = 0.667$$



Now, we have found cosine similarity between the test vector research and all other vectors and the word closest to research is :-

Theory with cosine similarity  
0.8164.

6) No, this is not a classification task as we have just represented our words as vectors and that doesn't inherently classify them.

Also, when we find similarity between 2 vectors, that doesn't automatically classify them into categories.

To use these vector embeddings for classification, we need to define what our classes are and the factors which will decide where our vector will belong.

Let us say we wish to classify each word as positive or negative then we can say that a vector whose norm is more than a specific value is positive or vice versa.

We can even use advanced classification techniques such as SVM's, Naive Bayes, Neural networks etc.

These methods will take our vector as an input and return the probability of being in one of the classes.

So we need to add a probabilistic training model or other training model to turn this into a classifier.

Q2) Rewriting the CV we obtain :-

	(data structures)	(eking ninjas)	(Prinly College)
doc1	0	0	0
doc2	1	0	0
doc3	1	1	0
doc4	0	1	0
doc5	1	0	0
doc6	0	0	1

  

	(College, London)	(structures, Algorithms)
doc1	0	0
doc2	0	1
doc3	0	1
doc4	0	0
doc5	0	1
doc6	1	0



By choosing tograms we obtain -

	(London, plectrum, guitar)	(Trinity, college, London)
doc1	0	0
doc2	0	0
doc3	0	0
doc4	0	0
doc5	0	0
doc6	1	1

	(data, structures, algorithms)	(plectrum, guitar, guitar)
doc1	1	0
doc2	1	0
doc3	0	0
doc4	1	0
doc5	0	0
doc6	0	1

	(college, London, plectrum)
doc1	0
doc2	0
doc3	0
doc4	0
doc5	0
doc6	1

The resume is a small document with not many words and by using higher order n-grams like 4-grams we lose words that occur highly in unigram but not in 4-gram like Java or python or

github.

So, we shouldn't use trigram and actually best result is delivered using unigram and bigram frequency.

We should prefer unigram over bigram as data like Java and Python is lost in bigrams, which is actually very important.

Hence, we should prefer unigram.



```
import nltk
import pandas
import pickle
from nltk.corpus import stopwords
from collections import Counter

stopwords_en = set(stopwords.words('english'))

# loading the document
resume = open('../assets/resume.txt', 'r')
document = resume.read().lower()
resume.close()

# tokenizing the document
tokenizer = nltk.RegexpTokenizer(r'\w+')
tokens = tokenizer.tokenize(document)

# removing stop words from tokens
tokens = [token for token in tokens if token not in stopwords_en and token.isalpha()]

# divide the document into n parts
documents = []
n = 6
k = len(tokens) // n
for i in range(n):
    documents.append(tokens[k * i: len(tokens) if i == n - 1 else (i + 1) * k])

def count_n_grams(data, n, start_token='<s>', end_token='<e>'):
    n_grams = {}
    for index, token in enumerate(data):
        n_gram = tuple(data[index: index + n])
        n_grams[n_gram] = n_grams.get(n_gram, 0) + 1
    return n_grams

bigrams = []
trigrams = []
for document in documents:
    bigrams.append(count_n_grams(tokens, 2))
    trigrams.append(count_n_grams(tokens, 3))

most_common_bi = set()
for bigram_freq in bigrams:
```



```

def count_n_grams(data, n, start_token='<s>', end_token='<e>'):
    n_grams = {}
    for index, token in enumerate(data):
        n_gram = tuple(data[index: index + n])
        n_grams[n_gram] = n_grams.get(n_gram, 0) + 1
    return n_grams

bigrams = []
trigrams = []
for document in documents:
    bigrams.append(count_n_grams(tokens, 2))
    trigrams.append(count_n_grams(tokens, 3))

most_common_bi = set()
for bigram_freq in bigrams:
    for word, freq in Counter(bigram_freq).most_common(5):
        most_common_bi.add(word)

most_common_tri = set()
for bigram_freq in trigrams:
    for word, freq in Counter(bigram_freq).most_common(5):
        most_common_tri.add(word)

print(most_common_bi)
print(most_common_tri)

vectors_bi = {}
for token in most_common_bi:
    vector = [0] * 6
    for index, bigram_freq in enumerate(bigrams):
        vector[index] = int(token in bigram_freq)
    vectors_bi[token] = vector

print(vectors_bi)

```



Q3) a) We can use wordnet to find all inflections of a word that are nouns and then see the different synonyms of a word and see these definitions eg. Sample code:-

Code for finding synonyms (Python)

```
from
import nltk.corpus import wordnet as wn
```

```
word = 'dog'
```

```
synsets = wn.synsets(word, pos='n')
```

```
>> print(synsets[0].definition)
```

```
A member of genus Canis (.....)
```

```
for synset in synsets:
    print(synset.definition())
```

output

A member of genus Canis...

A dull unattractive girl or woman

Disfornal Team for Man

Someone who is normally...

(We get various synonyms)

a) We can even incorporate WordNet to identify Lemmas and use it in Lemmatization process.

b) Using WordNet we will find the closest synonyms and also consider closest synonyms. We will reduce all words into it's lowest common hyperonym and consider those words,

So for dog and cat we consider → Carnivore.

The ~~more~~ number of synonyms will ~~increase~~ decrease but frequency will increase. It will lead to more generalization but loss of information.