

Assignment #5

CSCI 642 Natural Language Processing

Yeshe Reddipalli
15896337

Vector Representation of the model

```
test.py x
Courses > CSCI 642 Natural Language Processing > Assignment 5 > test.py > ...
1 import requests
2 from gensim import corpora, models
3 import numpy as np
4 from nltk.corpus import stopwords
5 from nltk.tokenize import word_tokenize
6
7 url = 'http://www.columbia.edu/~fdc/sample.html'
8 response = requests.get(url)
9 text = response.text
10
11 sentences = text.split('.')
12
13 stop_words = set(['a', 'an', 'the', 'is', 'are', 'was', 'were', 'in', 'on', 'at', 'to', 'of', 'and'])
14 sentences_preprocessed = []
15 for sentence in sentences:
16     words = sentence.strip().lower().split()
17     words_filtered = [word for word in words if word not in stop_words]
18     sentences_preprocessed.append(words_filtered)
19
20 dictionary = corpora.Dictionary(sentences_preprocessed)
21
22 LSA Vectors of the model:
23 Sentence 1: [ 0.535 -0.183 0.287]
24 Sentence 2: [ 0.01 -0.008 0.003]
25 Sentence 3: [ 0.004 -0.001 0.002]
26 Sentence 4: [ 0.004 -0.001 0.002]
27 Sentence 5: [ 0.274 0.007 -0.003]
28 Sentence 6: [ 0.257 0.045 -0.025]
29 Sentence 7: [ 0.257 0.045 -0.025]
30 Sentence 8: [ 0.57 -0.298 0.404]
31 Sentence 9: [ 0.03 -0.665 -0.659]
32 Sentence 10: []
33 Sentence 11: []
34 Sentence 12: [ 0.218 -0.271 0.155]
35 Sentence 13: [ 0.002 -0.012 0.008]
36 Sentence 14: [ 0.015 -0.073 0.099]
37 Sentence 15: [ 0.291 -0.781 0.76 ]
38 Sentence 16: [ 0.465 -2.456 2.361]
39 Sentence 17: [ 0.201 -0.362 0.489]
40 Sentence 18: [ 0.017 -0.17 0.248]
41 Sentence 19: [ 0.002 -0.012 0.008]
```

Top 3 topics

```
test.py x B5.py
Courses > CSCI 642 Natural Language Processing > Assignment 5 > test.py > ...
11 sentences = text.split('.')
12
13 stop_words = set(['a', 'an', 'the', 'is', 'are', 'was', 'were', 'in', 'on', 'at', 'to', 'of', 'and'])
14 sentences_preprocessed = []
15 for sentence in sentences:
16     words = sentence.strip().lower().split()
17     words_filtered = [word for word in words if word not in stop_words]
18     sentences_preprocessed.append(words_filtered)
19
20 dictionary = corpora.Dictionary(sentences_preprocessed)
21
22 bow_corpus = [dictionary.doc2bow(sentence) for sentence in sentences_preprocessed]
23
24 # Train the LSA model using the bag-of-words representation of the preprocessed sentences
25 lsa_model = models.LsiModel(bow_corpus, num_topics=3, id2word=dictionary)
26
27 # Get the vector representation of the preprocessed sentences based on the LSA model
28 lsa_vectors = []
29 for sentence in sentences_preprocessed:
30     lsa_vectors.append(lsa_model[dict(zip(dictionary.get_vocab(), sentence))])
31
32 Top 3 Topics of the model :
33 (0, '0.623*tdcell' + 0.335*ctr' + 0.251*table' + 0.239*-->' + 0.239*cl--' + 0.230*thpheading' + 0.204*cp' + 0.160*start' + 0.159*row' + 0.128*with')
34 (1, '-0.641*cspan' + -0.403*you' + -0.210*cp' + -0.207*your' + -0.187*web' + -0.150*can' + -0.141*lixa' + -0.136*if' + 0.123*tdcell' + -0.116*page')
35 (2, '-0.656*cspan' + 0.323*you' + 0.243*web' + 0.229*your' + 0.227*lixa' + 0.184*can' + 0.141*page' + 0.122*if' + 0.102*cp' + 0.093*have')
```

Code:

```
import requests
from gensim import corpora, models
import numpy as np
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

url = 'http://www.columbia.edu/~fdc/sample.html'
response = requests.get(url)
text = response.text

sentences = text.split('.')

stop_words = set(['a', 'an', 'the', 'is', 'are', 'was', 'were', 'in', 'on',
                  'at', 'to', 'of', 'and'])
sentences_preprocessed = []
for sentence in sentences:
    words = sentence.strip().lower().split()
    words_filtered = [word for word in words if word not in stop_words]
    sentences_preprocessed.append(words_filtered)

dictionary = corpora.Dictionary(sentences_preprocessed)

bow_corpus = [dictionary.doc2bow(sentence) for sentence in
               sentences_preprocessed]

# Train the LSA model using the bag-of-words representation of the
preprocessed sentences
lsa_model = models.LsiModel(bow_corpus, num_topics=3, id2word=dictionary)

# Get the vector representation of the preprocessed sentences based on the LSA
model
lsa_vectors = []
for sentence in sentences_preprocessed:
    lsa_vector = np.array([tup[1] for tup in
                           lsa_model[dictionary.doc2bow(sentence)]]
                           lsa_vectors.append(lsa_vector)

# Print the vector representation of the preprocessed sentences
print('LSA Vectors of the model:')
for i, vector in enumerate(lsa_vectors):
    print(f'Sentence {i+1}: {vector.round(3)}')

# Get the top 3 topics based on the LSA model
lsa_topics = lsa_model.print_topics(num_topics=3)
print('\nTop 3 Topics of the model :')
for topic in lsa_topics:
    print(topic)
```