Discussion Forums

# Week 4

| SUBFORUMS |
| --- |
| **All** |
| Assignment: Assignment 4 Submission |

← Week 4

## Week 4 Notebook Provided Here 📌                    ⌄

Uwe F Mayer   Mentor   Week 4 · 2 months ago

There isn't a provided one. I assembled a workbook when I took the class from the material presented. Here's the code, each block is a cell.

```python
1   import re
2   import pandas as pd
3   import numpy as np
4   import nltk
5   from nltk.corpus import wordnet as wn
6
7   # Use path length in wordnet to find word similarity
8   # find sense of words via synonym set
9   # n=noun, 01=synonym set for first meaning of the word
10  deer = wn.synset('deer.n.01')
11  deer
12
13  elk = wn.synset('elk.n.01')
14  deer.path_similarity(elk)
15
16  horse = wn.synset('horse.n.01')
17  deer.path_similarity(horse)
18
19  # Use an information criteria to find word similarity
20  from nltk.corpus import wordnet_ic
21  brown_ic = wordnet_ic.ic('ic-brown.dat')
22  deer.lin_similarity(elk, brown_ic)
23
24  deer.lin_similarity(horse, brown_ic)
25
26  # Use NLTK Collocation and Association Measures
27  from nltk.collocations import *
28  # load some text for examples
29  from nltk.book import *
30  # text1 is the book "Moby Dick"
31  # extract just the words without numbers and sentence marks and
        make them lower case
32  text = [w.lower() for w in list(text1) if w.isalpha()]
33
34  bigram_measures = nltk.collocations.BigramAssocMeasures()
35  finder = BigramCollocationFinder.from_words(text)
36  finder.nbest(bigram_measures.pmi,10)
37
38  # find all the bigrams with occurrence of at least 10, this
        modifies our "finder" object
39  finder.apply_freq_filter(10)
40  finder.nbest(bigram_measures.pmi,10)
41
42  # Working with Latent Dirichlet Allocation (LDA) in Python
43  # Several packages available, such as gensim and lda. Text needs
        to be
44  # preprocessed: tokenizing, normalizing such as lower-casing,
        stopword
45  # removal, stemming, and then transforming into a (sparse)
        matrix for
46  # word (bigram, etc) occurences.
47  # generate a set of preprocessed documents
48  from nltk.stem.porter import PorterStemmer
49  from nltk.corpus import stopwords
50  from nltk.book import *
51
52  len(stopwords.words('english'))
53
54  stopwords.words('english')
55
56  # extract just the stemmed words without numbers and sentence
        marks and make them lower case
57  p_stemmer = PorterStemmer()
58  sw = stopwords.words('english')
59  doc1 = [p_stemmer.stem(w.lower()) for w in list(text1) if w
        .isalpha() and not w.lower() in sw]
60  doc2 = [p_stemmer.stem(w.lower()) for w in list(text2) if w
        .isalpha() and not w.lower() in sw]
61  doc3 = [p_stemmer.stem(w.lower()) for w in list(text3) if w
        .isalpha() and not w.lower() in sw]
62  doc4 = [p_stemmer.stem(w.lower()) for w in list(text4) if w
        .isalpha() and not w.lower() in sw]
63  doc5 = [p_stemmer.stem(w.lower()) for w in list(text5) if w
        .isalpha() and not w.lower() in sw]
64  doc_set = [doc1, doc2, doc3, doc4, doc5]
65
66  # under Windows this generates a warning
67  import gensim
68  from gensim import corpora, models
69
```

```
70   dictionary = corpora.Dictionary(doc_set)
71   dictionary
72
73   # transform each document into a bag of words
74   corpus = [dictionary.doc2bow((doc)) for doc in doc_set]
75
76   # The corpus contains the 5 documents
77   # each document is a list of indexed features and occurrence
         count (freq)
78   print(type(corpus))
79   print(type(corpus[0]))
80   print(type(corpus[0][0]))
81   print(corpus[0][::2000])
82
83   # let's try 4 topics for our 5 documents
84   # 50 passes takes quite a while, let's try less
85   ldamodel = gensim.models.ldamodel.LdaModel(corpus, num_topics=4,
         id2word=dictionary, passes=10)
86
87   print(ldamodel.print_topics(num_topics=4, num_words=10))
```

⇧ 11 Upvotes        💬 Reply        Follow this discussion

---

**Earliest**              **Top**                 **Most Recent**

Kedar Joshi · 5 days ago

This certainly was a great help! Thanks Uwe and thank for your insights overall on various discussion threads throughout the course.

⇧ 0 Upvotes        💬 Reply

OP   Oscar Rene Chamberlain Pravia · a month ago

Excellent contribution!!

⇧ 0 Upvotes        💬 Reply

Kakoli · a month ago

This is a great help.

⇧ 0 Upvotes        💬 Reply

‹  [1]  ›

HW   | Reply
     |
     |

Reply