

```
1 #%% md
2 ## NLP with DL Assignment 1:
3 # Exploring Word Vectors
4 ### <font color='red'> Due 12:00am, Sun April 27 </font>
5
6
7
8 Before you start, make sure you read the README.txt
  in the same directory as this notebook for
  important setup information. A lot of code is
  provided in this notebook, and we highly encourage
  you to read and understand it as part of the
  learning :)
9
10 If you aren't super familiar with Python, Numpy, or
    Matplotlib, we recommend you check out a Stanford'
    s [tutorial](https://cs231n.github.io/python-numpy-
    tutorial/) about Python/Numpy.
11
12 **Assignment Notes:** Please make sure to save the
    notebook as you go along. Submission Instructions
    are located at the bottom of the notebook.
13 #%%
14 # All Import Statements Defined Here
15 # Note: Do not add to this list.
16 #
17
18 import sys
19 assert sys.version_info[0]==3
20 assert sys.version_info[1] >= 5
21
22 from platform import python_version
23 assert int(python_version().split(".")[1]) >= 5, "
  Please upgrade your Python version following the
  instructions in \
24     the README.txt file found in the same directory
  as this notebook. Your Python version is " +
  python_version()
25
26 from gensim.models import KeyedVectors
```

```
27 from gensim.test.utils import datapath
28 import pprint
29 import matplotlib.pyplot as plt
30 plt.rcParams['figure.figsize'] = [10, 5]
31
32 import nltk
33 nltk.download('reuters') #to specify download
   location, optionally add the argument: download_dir
   ='./specify/desired/path/'
34 from nltk.corpus import reuters
35
36 import numpy as np
37 import random
38 import scipy as sp
39 from sklearn.decomposition import TruncatedSVD
40 from sklearn.decomposition import PCA
41
42 START_TOKEN = '<START>'
43 END_TOKEN = '<END>'
44
45 np.random.seed(0)
46 random.seed(0)
47 # -----
48 #%% md
49 ## Word Vectors
50
51 Word Vectors are often used as a fundamental
   component for downstream NLP tasks, e.g. question
   answering, text generation, translation, etc., so
   it is important to build some intuitions as to
   their strengths and weaknesses. Here, you will
   explore two types of word vectors: those derived
   from *co-occurrence matrices*, and those derived
   via *GloVe*.
52
53 **Note on Terminology:** The terms "word vectors"
   and "word embeddings" are often used
   interchangeably. The term "embedding" refers to the
   fact that we are encoding aspects of a word's
   meaning in a lower dimensional space. As [Wikipedia
   ](https://en.wikipedia.org/wiki/Word\_embedding)
```

```

53 states, "*conceptually it involves a mathematical
embedding from a space with one dimension per word
to a continuous vector space with a much lower
dimension*".


---


54 #%% md
55 ## Part 1: Count-Based Word Vectors
56
57 Most word vector models start from the following
idea:
58
59 *You shall know a word by the company it keeps ([  

  Firth, J. R. 1957:11](https://en.wikipedia.org/wiki/John\_Rupert\_Firth))*  

60
61 Many word vector implementations are driven by the
idea that similar words, i.e., (near) synonyms,
will be used in similar contexts. As a result,
similar words will often be spoken or written along
with a shared subset of words, i.e., contexts. By
examining these contexts, we can try to develop
embeddings for our words. With this intuition in
mind, many "old school" approaches to constructing
word vectors relied on word counts. Here we
elaborate upon one of those strategies, *co-
occurrence matrices* (for more information, see
Word Embedding lecture).


---


62 #%% md
63 ### Co-Occurrence
64
65 A co-occurrence matrix counts how often things co-
occur in some environment. Given some word  $w_i$ 
occurring in the document, we consider the *context
window* surrounding  $w_i$ . Supposing our fixed
window size is  $n$ , then this is the  $n$  preceding
and  $n$  subsequent words in that document, i.e.
words  $w_{i-n} \dots w_{i-1}$  and  $w_{i+1} \dots w_{i+n}$ . We build a *co-occurrence matrix*  $M$ ,
which is a symmetric word-by-word matrix in which
 $M_{ij}$  is the number of times  $w_j$  appears
inside  $w_i$ 's window among all documents.
66

```

```

67 **Example: Co-Occurrence with Fixed Window of n=1
**:
68
69 Document 1: "all that glitters is not gold"
70
71 Document 2: "all is well that ends well"
72
73
74 | * | `<START>` | all | that | glitters |
  is | not | gold | well | ends | `<END>` |
75 |-----|-----|-----|-----|-----|-----
  |-----|-----|-----|-----|-----|
76 | `<START>` | 0 | 2 | 0 | 0 | 0 | 0
  | 0 | 0 | 0 | 0 | 0 | 0 |
77 | all | 2 | 0 | 1 | 0 | 0 | 1
  | 0 | 0 | 0 | 0 | 0 | 0 |
78 | that | 0 | 1 | 0 | 1 | 0 | 0
  | 0 | 0 | 1 | 1 | 0 | 0 |
79 | glitters | 0 | 0 | 1 | 0 | 0 | 1
  | 0 | 0 | 0 | 0 | 0 | 0 |
80 | is | 0 | 1 | 0 | 1 | 0 | 0
  | 1 | 0 | 1 | 0 | 0 | 0 |
81 | not | 0 | 0 | 0 | 0 | 0 | 1
  | 0 | 1 | 0 | 0 | 0 | 0 |
82 | gold | 0 | 0 | 0 | 0 | 0 | 0
  | 1 | 0 | 0 | 0 | 1 | 0 |
83 | well | 0 | 0 | 1 | 0 | 0 | 1
  | 0 | 0 | 0 | 1 | 1 | 0 |
84 | ends | 0 | 0 | 1 | 0 | 0 | 0
  | 0 | 0 | 1 | 0 | 0 | 0 |
85 | `<END>` | 0 | 0 | 0 | 0 | 0 | 0
  | 0 | 1 | 1 | 0 | 0 | 0 |
86
87 **Note:** In NLP, we often add `<START>` and `<END>` tokens to represent the beginning and end of sentences, paragraphs or documents. In this case we imagine `<START>` and `<END>` tokens encapsulating each document, e.g., "`<START>` All that glitters is not gold `<END>`", and include these tokens in our co-occurrence counts.
88

```

```

89 The rows (or columns) of this matrix provide one
    type of word vectors (those based on word-word co-
    occurrence), but the vectors will be large in
    general (linear in the number of distinct words in
    a corpus). Thus, our next step is to run *
    dimensionality reduction*. In particular, we will
    run *SVD (Singular Value Decomposition)*, which is
    a kind of generalized *PCA (Principal Components
    Analysis)* to select the top $k$ principal
    components. Here's a visualization of
    dimensionality reduction with SVD. In this picture
    our co-occurrence matrix is $A$ with $n$ rows
    corresponding to $n$ words. We obtain a full
    matrix decomposition, with the singular values
    ordered in the diagonal $$ matrix, and our new,
    shorter length-$k$ word vectors in $U_k$.

90
91 ! [Picture of an SVD](./imgs/svd.png "SVD")
92
93 This reduced-dimensionality co-occurrence
    representation preserves semantic relationships
    between words, e.g. *doctor* and *hospital* will
    be closer than *doctor* and *dog*.
94
95 **Notes:** If you can barely remember what an
    eigenvalue is, here's [a slow, friendly
    introduction to SVD](https://davetang.org/file/Singular\_Value\_Decomposition\_Tutorial.pdf). In
    practice, it is challenging to apply full SVD to
    large corpora because of the memory needed to
    perform PCA or SVD. However, if you only want the
    top $k$ vector components for relatively small $k$-
    – known as [Truncated SVD](https://en.wikipedia.org/wiki/Singular\_value\_decomposition#Truncated\_SVD) – then there are reasonably
    scalable techniques to compute those iteratively.
96 #%% md
97 ### Plotting Co-Occurrence Word Embeddings
98
99 Here, we will be using the Reuters (business and
    financial news) corpus. If you haven't run the

```

```

99 import cell at the top of this page, please run it
    now (click it and press SHIFT-RETURN). The corpus
    consists of 10,788 news documents totaling 1.3
    million words. These documents span 90 categories
    and are split into train and test. For more
    details, please see https://www.nltk.org/book/ch02.html. We provide a `read_corpus` function below
    that pulls out only articles from the "gold" (i.e.
    . news articles about gold, mining, etc.) category
    . The function also adds `<START>` and `<END>`
    tokens to each of the documents, and lowercases
    words. You do **not** have to perform any other
    kind of pre-processing.
100 #%%
101 def read_corpus(category="gold"):
102     """ Read files from the specified Reuter's
        category.
103     Params:
104         category (string): category name
105     Return:
106         list of lists, with words from each of
        the processed files
107     """
108     files = reuters.fileids(category)
109     return [[START_TOKEN] + [w.lower() for w in
        list(reuters.words(f))] + [END_TOKEN] for f in
        files]
110
111 #%% md
112 Let's have a look what these documents are like...
113 #%%
114 reuters_corpus = read_corpus()
115 pprint.pprint(reuters_corpus[:3], compact=True,
    width=100)
116 #%% md
117 ### Question 1.1: Implement `distinct_words` [code
    ]
118
119 Write a method to work out the distinct words (
    word types) that occur in the corpus. You can do
    this with `for` loops, but it's more efficient to

```

```

119 do it with Python list comprehensions. In
particular, [this](https://coderwall.com/p/rcmaea/
flatten-a-list-of-lists-in-one-line-in-python) may
be useful to flatten a list of lists. If you're
not familiar with Python list comprehensions in
general, here's [more information](https://python-
3-patterns-idioms-test.readthedocs.io/en/latest/
Comprehensions.html).
120
121 Your returned `corpus_words` should be sorted. You
can use python's `sorted` function for this.
122
123 You may find it useful to use [Python sets](https
://www.w3schools.com/python/python_sets.asp) to
remove duplicate words.
124 #%%
125 def distinct_words(corpus):
126     """ Determine a list of distinct words for the
corpus.
127         Params:
128             corpus (list of list of strings):
corpus of documents
129         Return:
130             corpus_words (list of strings): sorted
list of distinct words across the corpus
131             n_corpus_words (integer): number of
distinct words across the corpus
132         """
133     corpus_words = []
134     n_corpus_words = -1
135
136     ### SOLUTION BEGIN
137     flattened_list = [y for x in corpus for y in x
]
138     corpus_words = sorted(set(flattened_list))
139     n_corpus_words = len(corpus_words)
140     ### SOLUTION END
141
142     return corpus_words, n_corpus_words
143 # -----

```

```

145 # Run this sanity check
146 # Note that this not an exhaustive check for
   correctness.
147 #
148
149 # Define toy corpus
150 test_corpus = ["{} All that glitters isn't gold
                  {}".format(START_TOKEN, END_TOKEN).split(" "),
                  "{} All's well that ends well {}".format(
                      START_TOKEN, END_TOKEN).split(" ")]
151
152 test_corpus_words, num_corpus_words =
   distinct_words(test_corpus)
153
154
155 # Correct answers
156 ans_test_corpus_words = sorted([START_TOKEN, "All",
                                   "ends", "that", "gold", "All's",
                                   "glitters", "isn't", "well",
                                   END_TOKEN])
157 ans_num_corpus_words = len(ans_test_corpus_words)
158
159 # Test correct number of words
160 assert(num_corpus_words == ans_num_corpus_words),
   "Incorrect number of distinct words. Correct: {}.
   Yours: {}".format(ans_num_corpus_words,
                     num_corpus_words)
161
162 # Test correct words
163 assert (test_corpus_words == ans_test_corpus_words),
   "Incorrect corpus_words.\nCorrect: {}\nYours
   : {}".format(str(ans_test_corpus_words), str(
     test_corpus_words))
164
165 # Print Success
166 print ("-" * 80)
167 print("Passed All Tests!")
168 print ("-" * 80)
169 #%% md
170 ### Question 1.2: Implement `compute_co_occurrence_matrix` [code]
171

```

```

172 Write a method that constructs a co-occurrence
    matrix for a certain window-size $n$ (with a
    default of 4), considering words $n$ before and
    $n$ after the word in the center of the window.
    Here, we start to use `numpy (np)` to represent
    vectors, matrices, and tensors.
173
174 #%%
175 def compute_co_occurrence_matrix(corpus,
176     window_size=4):
177     """ Compute co-occurrence matrix for the given
        corpus and window_size (default of 4).
178
179         Note: Each word in a document should be at
            the center of a window. Words near edges will
            have a smaller
180             number of co-occurring words.
181
182             For example, if we take the document
            "<START> All that glitters is not gold <END>""
            with window size of 4,
183                 "All" will co-occur with "<START
            >", "that", "glitters", "is", and "not".
184
185             Params:
186                 corpus (list of list of strings):
            corpus of documents
187                 window_size (int): size of context
            window
188
189             Return:
190                 M (a symmetric numpy matrix of shape (
            number of unique words in the corpus , number of
            unique words in the corpus)):
191                 Co-occurrence matrix of word counts
192
193                 .
194
195                 The ordering of the words in the
            rows/columns should be the same as the ordering of
            the words given by the distinct_words function.
196
197                 word2ind (dict): dictionary that maps
            word to index (i.e. row/column number) for matrix
            M.

```

```

192     """
193     words, n_words = distinct_words(corpus)
194     M = None
195     word2ind = {}
196
197     ### SOLUTION BEGIN
198     # create initial M in size nXn
199     M = np.zeros((n_words, n_words))
200     word2ind = dict(zip(words, range(n_words)))
201
202     for document in corpus:
203         for i, center_word in enumerate(document):
204             center_index = word2ind[center_word]
205             left_bound = max(0, i - window_size)
206             right_bound = min(len(document), i +
207                               window_size + 1)
208
209             for j in range(left_bound, right_bound):
210                 if i != j:
211                     neighbor_word = document[j]
212                     neighbor_index = word2ind[
213                         neighbor_word]
214                     M[center_index, neighbor_index]
215                     ] += 1
216
217     # ### SOLUTION END
218     #
219     return M, word2ind
220 # -----
221 # Run this sanity check
222 # Note that this is not an exhaustive check for
223 # correctness.
224 # -----
225 # Define toy corpus and get student's co-
226 # occurrence matrix
227 test_corpus = ["{} All that glitters isn't gold
228                 {}".format(START_TOKEN, END_TOKEN).split(" "),

```

```

225 "{} All's well that ends well {}".format(
    START_TOKEN, END_TOKEN).split(" ")]
226 M_test, word2ind_test =
    compute_co_occurrence_matrix(test_corpus,
        window_size=1)
227
228 # Correct M and word2ind
229 M_test_ans = np.array(
230     [[0., 0., 0., 0., 0., 0., 1., 0., 0., 1.],
231      [0., 0., 1., 1., 0., 0., 0., 0., 0., 0.],
232      [0., 1., 0., 0., 0., 0., 0., 0., 1., 0.],
233      [0., 1., 0., 0., 0., 0., 0., 0., 0., 1.],
234      [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.],
235      [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.],
236      [1., 0., 0., 0., 0., 0., 0., 1., 0., 0.],
237      [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.],
238      [0., 0., 1., 0., 1., 1., 0., 0., 0., 1.],
239      [1., 0., 0., 1., 1., 0., 0., 1., 0., 0.]])
240 )
241 ans_test_corpus_words = sorted([START_TOKEN, "All",
        "ends", "that", "gold", "All's", "glitters",
        "isn't", "well", END_TOKEN])
242 word2ind_ans = dict(zip(ans_test_corpus_words,
    range(len(ans_test_corpus_words))))
243
244 # Test correct word2ind
245 assert (word2ind_ans == word2ind_test), "Your
    word2ind is incorrect:\nCorrect: {}\nYours: {}".
    format(word2ind_ans, word2ind_test)
246
247 # Test correct M shape
248 assert (M_test.shape == M_test_ans.shape), "M
    matrix has incorrect shape.\nCorrect: {}\nYours
    : {}".format(M_test.shape, M_test_ans.shape)
249
250 # Test correct M values
251 for w1 in word2ind_ans.keys():
252     idx1 = word2ind_ans[w1]
253     for w2 in word2ind_ans.keys():
254         idx2 = word2ind_ans[w2]
255         student = M_test[idx1, idx2]

```

```

256     correct = M_test_ans[idx1, idx2]
257     if student != correct:
258         print("Correct M:")
259         print(M_test_ans)
260         print("Your M: ")
261         print(M_test)
262         raise AssertionError("Incorrect count
263             at index ({}, {}) = ({}, {}) in matrix M. Yours has
264             {} but should have {}.".format(idx1, idx2, w1, w2
265             , student, correct))
266
267 # Print Success
268 print ("-" * 80)
269 print("Passed All Tests!")
270 print ("-" * 80)
271 #### Question 1.3: Implement `reduce_to_k_dim` [code]
272
273 **Note:** All of numpy, scipy, and scikit-learn (`sklearn`) provide *some* implementation of SVD, but only scipy and sklearn provide an implementation of Truncated SVD, and only sklearn provides an efficient randomized algorithm for calculating large-scale Truncated SVD. So please use [sklearn.decomposition.TruncatedSVD](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html).
274 %%#
275 def reduce_to_k_dim(M, k=2):
276     """Reduce a co-occurrence count matrix of
277     dimensionality (num_corpus_words, num_corpus_words)
278         to a matrix of dimensionality (
279         num_corpus_words, k) using the following SVD

```

```

277 function from Scikit-Learn:
278         - http://scikit-learn.org/stable/
279             modules/generated/sklearn.decomposition.
280                 TruncatedSVD.html
281
282             Params:
283                 M (numpy matrix of shape (number of
284                     unique words in the corpus , number of unique
285                     words in the corpus)): co-occurrence matrix of word
286                     counts
287                     k (int): embedding size of each word
288                     after dimension reduction
289             Return:
290                 M_reduced (numpy matrix of shape (
291                     number of corpus words, k)): matrix of k-
292                     dimensioal word embeddings.
293                     In terms of the SVD from math
294                     class, this actually returns U * S
295                     """
296                     n_iters = 10      # Use this parameter in your
297                     call to `TruncatedSVD`
298                     M_reduced = None
299                     print("Running Truncated SVD over %i words..." %
300                         (M.shape[0]))
301
302                     #### SOLUTION BEGIN
303                     svd = TruncatedSVD(n_components=k, n_iter=
304                         n_iters)
305                     M_reduced = svd.fit_transform(M)
306                     #### SOLUTION END
307
308                     print("Done.")
309                     return M_reduced
310
311 #%%
312 # -----
313 # Run this sanity check
314 # Note that this is not an exhaustive check for
315 # correctness
316 # In fact we only check that your M_reduced has
317 # the right dimensions.
318 # -----

```

```

304
305 # Define toy corpus and run student code
306 test_corpus = ["{} All that glitters isn't gold
                  {}".format(START_TOKEN, END_TOKEN).split(" "),
                  "{} All's well that ends well {}".format(
                      START_TOKEN, END_TOKEN).split(" ")]
307 M_test, word2ind_test =
            compute_co_occurrence_matrix(test_corpus,
                                          window_size=1)
308 M_test_reduced = reduce_to_k_dim(M_test, k=2)
309
310 # Test proper dimensions
311 assert (M_test_reduced.shape[0] == 10), "M_reduced
                                             has {} rows; should have {}".format(
                                             M_test_reduced.shape[0], 10)
312 assert (M_test_reduced.shape[1] == 2), "M_reduced
                                             has {} columns; should have {}".format(
                                             M_test_reduced.shape[1], 2)
313
314 # Print Success
315 print ("-" * 80)
316 print("Passed All Tests!")
317 print ("-" * 80)
318 #%% md
319 ### Question 1.4: Implement `plot_embeddings` [code]
320
321 Here you will write a function to plot a set of 2D
      vectors in 2D space. For graphs, we will use
      Matplotlib (`plt`).
322
323 For this example, you may find it useful to adapt
      [this code](http://web.archive.org/web/
      20190924160434/https://www.pythonguides.com/2018
      /05/08/matplotlib-scatter-plot-annotate-set-text-
      at-label-each-point/). In the future, a good way
      to make a plot is to look at [the Matplotlib
      gallery](https://matplotlib.org/gallery/index.html
      ), find a plot that looks somewhat like what you
      want, and adapt the code they give.
324 #%%

```

```

325
326 def plot_embeddings(M_reduced, word2ind, words):
327     """ Plot in a scatterplot the embeddings of
328         the words specified in the list "words".
329         NOTE: do not plot all the words listed in
329         M_reduced / word2ind.
330
331             Include a label next to each point.
330
331     Params:
332         M_reduced (numpy matrix of shape (
333             number of unique words in the corpus , 2)): matrix
334             of 2-dimensioal word embeddings
333         word2ind (dict): dictionary that maps
334             word to indices for matrix M
334         words (list of strings): words whose
334             embeddings we want to visualize
335     """
336
337     ### SOLUTION BEGIN
338     plt.figure(figsize=(10, 10))
339     for word in words:
340         if word in word2ind:
341             index = word2ind[word]
342             embedding = M_reduced[index]
343             plt.scatter(embedding[0], embedding[1]
344             )
344             plt.annotate(word, xy=(embedding[0],
344             embedding[1]), xytext=(5, 2), textcoords='offset
344             points', ha='right', va='bottom')
345     plt.show()
346     ### SOLUTION END
347 #%%
348 # -----
349 # Run this sanity check
350 # Note that this is not an exhaustive check for
350 # correctness.
351 # The plot produced should look like the "test
351 # solution plot" depicted below.
352 # -----
353
354 print ("-" * 80)

```

```

355 print ("Outputted Plot:")
356
357 M_reduced_plot_test = np.array([[1, 1], [-1, -1]
358 [1, -1], [-1, 1], [0, 0]])
358 word2ind_plot_test = {'test1': 0, 'test2': 1,
359 'test3': 2, 'test4': 3, 'test5': 4}
359 words = ['test1', 'test2', 'test3', 'test4',
360 'test5']
360 plot_embeddings(M_reduced_plot_test,
361 word2ind_plot_test, words)
361
362 print ("-" * 80)
363 #%% md
364 ### Question 1.5: Co-Occurrence Plot Analysis [
364 written]
365
366 Now we will put together all the parts you have
366 written! We will compute the co-occurrence matrix
366 with fixed window of 4 (the default window size),
366 over the Reuters "gold" corpus. Then we will use
366 TruncatedSVD to compute 2-dimensional embeddings
366 of each word. TruncatedSVD returns  $U \times S$ , so we
366 need to normalize the returned vectors, so that
366 all the vectors will appear around the unit circle
366 (therefore closeness is directional closeness
366 ). **Note**: The line of code below that does the
366 normalizing uses the NumPy concept of *
366 broadcasting*. If you don't know about
366 broadcasting, check out
367 [Computation on Arrays: Broadcasting by Jake
367 VanderPlas](https://jakevdp.github.io/
367 PythonDataScienceHandbook/02.05-computation-on-
367 arrays-broadcasting.html).
368
369 Run the below cell to produce the plot. It'll
369 probably take a few seconds to run.
370 #%%
371 # -----
372 # Run This Cell to Produce Your Plot
373 # -----
374 reuters_corpus = read_corpus()

```

```

375 M_co_occurrence, word2ind_co_occurrence =
    compute_co_occurrence_matrix(reuters_corpus)
376 M_reduced_co_occurrence = reduce_to_k_dim(
    M_co_occurrence, k=2)
377
378 # Rescale (normalize) the rows to make them each
   of unit-length
379 M_lengths = np.linalg.norm(M_reduced_co_occurrence
   , axis=1)
380 M_normalized = M_reduced_co_occurrence / M_lengths
   [:, np.newaxis] # broadcasting
381
382 words = ['value', 'gold', 'platinum', 'reserves',
   'silver', 'metals', 'copper', 'belgium',
   'australia', 'china', 'grammes', "mine"]
383
384 plot_embeddings(M_normalized,
   word2ind_co_occurrence, words)
385 #%% md
386 **Verify that your figure matches "question_1.5.
   png" in the assignment zip. If not, use that
   figure to answer the next two questions.**
387 #%% md
388 a. Find at least two groups of words that cluster
   together in 2-dimensional embedding space. Give an
   explanation for each cluster you observe.
389 #%% md
390 ### SOLUTION BEGIN
391 Groups of words that cluster together in 2-
   dimensional embedding space:
392 1. The group of 'copper', 'platinum' & 'silver'
   which are all types of metals cluster together.
   However, silver appears slightly farther from
   copper and platinum, which makes sense given the
   broader semantic silver might have compared to the
   other two.
393 2. The group of 'belgium' & 'australia' form a
   tight cluster, which aligns with their semantic
   similarity as country names.
394
395 ### SOLUTION END

```

```
396 _____  
397 #%% md  
398 b. What doesn't cluster together that you might  
think should have? Describe at least two examples.  
399 #%% md  
400 ### SOLUTION BEGIN  
401  
402 1. I would expect that 'metals' and 'gold' will  
also cluster with the group of 'copper', '  
platinum' & 'silver', due to their semantic  
relations.  
403 2. I would expect that also 'china' will be in the  
cluster of 'belgium' & 'australia' as a country  
name as well.  
404 ### SOLUTION END  
405 #%% md  
406 ## Part 2: Prediction-Based Word Vectors  
407  
408 As discussed in class, more recently prediction-  
based word vectors have demonstrated better  
performance, such as word2vec and GloVe (which  
also utilizes the benefit of counts). Here, we  
shall explore the embeddings produced by GloVe.  
Please revisit the class notes and lecture slides  
for more details on the word2vec and GloVe  
algorithms. If you're feeling adventurous,  
challenge yourself and try reading [GloVe's  
original paper](https://nlp.stanford.edu/pubs/glove.pdf).  
409  
410 Then run the following cells to load the GloVe  
vectors into memory. **Note**: If this is your  
first time to run these cells, i.e. download the  
embedding model, it will take a couple minutes to  
run. If you've run these cells before, rerunning  
them will load the model without redownloading it  
, which will take about 1 to 2 minutes.  
411 #%%  
412 def load_embedding_model():  
413     """ Load GloVe Vectors  
414         Return:
```

```

415      wv_from_bin: All 400000 embeddings,
416      each length 200
417      """
418      import gensim.downloader as api
419      wv_from_bin = api.load("glove-wiki-gigaword-
420      200")
421      print("Loaded vocab size %i" % len(list(
422      wv_from_bin.index_to_key)))
423      return wv_from_bin
424      #%%%
425      # -----
426      # Run Cell to Load Word Vectors
427      # Note: This will take a couple minutes
428      # -----
429      #%% md
430      ##### Note: If you are receiving a "reset by peer"
431      error, rerun the cell to restart the download. If
432      you run into an "attribute" error, you may need to
433      update to the most recent version of gensim and
434      numpy. You can upgrade them inline by uncommenting
435      and running the below cell:
436      #%%%
437      #!pip install gensim --upgrade
438      #!pip install numpy --upgrade
439      #%% md
440      ### Reducing dimensionality of Word Embeddings
441      Let's directly compare the GloVe embeddings to
442      those of the co-occurrence matrix. In order to
443      avoid running out of memory, we will work with a
444      sample of 10000 GloVe vectors instead.
445      Run the following cells to:
446
447      1. Put 10000 Glove vectors into a matrix M
448      2. Run `reduce_to_k_dim` (your Truncated SVD
449      function) to reduce the vectors from 200-
450      dimensional to 2-dimensional.
451      #%%
452      def get_matrix_of_vectors(wv_from_bin,
453      required_words):
454          """ Put the GloVe vectors into a matrix M.

```

```

442         Param:
443             wv_from_bin: KeyedVectors object; the
444                 400000 GloVe vectors loaded from file
445             Return:
446                 M: numpy matrix shape (num words, 200
447                     ) containing the vectors
448                 word2ind: dictionary mapping each word
449                     to its row number in M
450                 """
451
452     import random
453     words = list(wv_from_bin.index_to_key)
454     print("Shuffling words ...")
455     random.seed(225)
456     random.shuffle(words)
457     words = words[:10000]
458     print("Putting %i words into word2ind and" %
459           matrix M..." % len(words))
460     word2ind = {}
461     M = []
462     curInd = 0
463     for w in words:
464         try:
465             M.append(wv_from_bin.get_vector(w))
466             word2ind[w] = curInd
467             curInd += 1
468         except KeyError:
469             continue
470     for w in required_words:
471         if w in words:
472             continue
473         try:
474             M.append(wv_from_bin.get_vector(w))
475             word2ind[w] = curInd
476             curInd += 1
477         except KeyError:
478             continue
479     M = np.stack(M)
480     print("Done.")
481     return M, word2ind
482
483 # -----

```

```

478 -----
479 # Run Cell to Reduce 200-Dimensional Word
  Embeddings to k Dimensions
480 # Note: This should be quick to run
481 # -----
482 M, word2ind = get_matrix_of_vectors(wv_from_bin,
  words)
483 M_reduced = reduce_to_k_dim(M, k=2)
484
485 # Rescale (normalize) the rows to make them each
  of unit-length
486 M_lengths = np.linalg.norm(M_reduced, axis=1)
487 M_reduced_normalized = M_reduced / M_lengths[:, np
  .newaxis] # broadcasting
488 #%% md
489 **Note: If you are receiving out of memory issues
  on your local machine, try closing other
  applications to free more memory on your device.
  You may want to try restarting your machine so
  that you can free up extra memory. Then
  immediately run the jupyter notebook and see if
  you can load the word vectors properly. If you
  still have problems with loading the embeddings
  onto your local machine after this, please go to
  office hours or contact course staff.**
490 #%% md
491 ### Question 2.1: GloVe Plot Analysis [written]
492
493 Run the cell below to plot the 2D GloVe embeddings
  for `['value', 'gold', 'platinum', 'reserves',
  'silver', 'metals', 'copper', 'belgium', 'australia',
  'china', 'grammes', "mine"]`.
494 #%%
495 words = ['value', 'gold', 'platinum', 'reserves',
  'silver', 'metals', 'copper', 'belgium', 'australia',
  'china', 'grammes', "mine"]
496
497 plot_embeddings(M_reduced_normalized, word2ind,
  words)
498 #%% md

```

499 a. What is one way the plot is different from the one generated earlier from the co-occurrence matrix? What is one way it's similar?

500 #%% md

501 ### SOLUTION BEGIN

502 1. This GloVe plot different by clustering silver , metal, and gold together as expected given their strong semantic relationships.

503 2. This GloVe plot is similar to the previous one in that it clusters australia and belgium together , while china isn't cluster with them as might expect.

504 ### SOLUTION END

505 #%% md

506 b. What is a possible cause for the difference?

507 #%% md

508 ### SOLUTION BEGIN

509 Since GloVe leverages global word co-occurrence statistics across the entire corpus, the differences between this plot and the previous one may reflect broader contextual similarities, rather than just local window-based associations.

510 ### SOLUTION END

511 #%% md

512 ### Cosine Similarity

513 Now that we have word vectors, we need a way to quantify the similarity between individual words, according to these vectors. One such metric is cosine-similarity. We will be using this to find words that are "close" and "far" from one another.

514

515 We can think of n-dimensional vectors as points in n-dimensional space. If we take this perspective [L1](<http://mathworld.wolfram.com/L1-Norm.html>) and [L2](<http://mathworld.wolfram.com/L2-Norm.html>) Distances help quantify the amount of space "we must travel" to get between these two points. Another approach is to examine the angle between two vectors. From trigonometry we know that:

516

517
518
519 Instead of computing the actual angle, we can
 leave the similarity in terms of $\text{similarity} = \cos(\Theta)$. Formally the [Cosine Similarity](https://en.wikipedia.org/wiki/Cosine_similarity) s
 between two vectors p and q is defined as:
520
521
$$s = \frac{p \cdot q}{\|p\| \|q\|}, \text{ where } s \in [-1, 1]$$

522 #%% md
523 ### Question 2.2: Words with Multiple Meanings [code + written]
524 Polysemes and homonyms are words that have more than one meaning (see this [wiki page](https://en.wikipedia.org/wiki/Polysemy) to learn more about the difference between polysemes and homonyms). Find a word with *at least two different meanings* such that the top-10 most similar words (according to cosine similarity) contain related words from *both* meanings. For example, "leaves" has both "go_away" and "a_structure_of_a_plant" meaning in the top 10, and "scoop" has both "handed_waffle_cone" and "lowdown". You will probably need to try several polysemous or homonymic words before you find one.
525
526 Please state the word you discover and the multiple meanings that occur in the top 10. Why do you think many of the polysemous or homonymic words you tried didn't work (i.e. the top-10 most similar words only contain **one** of the meanings of the words)?
527
528 **Note**: You should use the `wv_from_bin.most_similar(word)` function to get the top 10 similar words. This function ranks all other words in the vocabulary with respect to their cosine similarity to the given word. For further assistance, please check the __[GenSim documentation](https://radimrehurek.com/gensim/)

```

```

528 models/keyedvectors.html#gensim.models.

 keyedvectors.FastTextKeyedVectors.most_similar)___.

529 #%%

530 ### SOLUTION BEGIN

531 wv_from_bin.most_similar("light")

532 ### SOLUTION END

533 #%% md

534 ### SOLUTION BEGIN

535 The word I discovered is 'light,' whose meanings

 include both 'brightness' and 'low weight' and

 more. This polysemy is reflected in its top

 similar words, such as 'bright' (number 1) and '

 heavy' (number 5).

536

537 In my opinion, many other polysemous or homonymic

 words didn't work because one of their meanings

 presented much more in the data and therefore in

 the embedding while the other is

 underrepresentation.

538 ### SOLUTION END

539 #%% md

540 ### Question 2.3: Synonyms & Antonyms [code +

 written]

541

542 When considering Cosine Similarity, it's often

 more convenient to think of Cosine Distance, which

 is simply $1 - \text{Cosine Similarity}$.

543

544 Find three words (w_1, w_2, w_3) where w_1 and

 w_2 are synonyms and w_1 and w_3 are

 antonyms, but Cosine Distance $(w_1, w_3) < \text{Cosine Distance } (w_1, w_2)$.

545

546 As an example, $w_1 = "happy"$ is closer to $w_3 = "sad"$ than to $w_2 = "cheerful"$. Please find a

 different example that satisfies the above. Once

 you have found your example, please give a

 possible explanation for why this counter-

 intuitive result may have happened.

547

548 You should use the the `wv_from_bin.distance(w1,

```

```

548 w2)` function here in order to compute the cosine
 distance between two words. Please see the __[

 GenSim documentation](https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors.FastTextKeyedVectors.distance)__ for

 further assistance.

549 #%%
550 ### SOLUTION BEGIN
551
552 w1 = "light"
553 w2 = "illuminate"
554 w3 = "heavy"
555
556 w1_w2_dist = wv_from_bin.distance(w1, w2)
557 w1_w3_dist = wv_from_bin.distance(w1, w3)
558
559 print("Synonyms {}, {} have cosine distance: {}".
 format(w1, w2, w1_w2_dist))
560 print("Antonyms {}, {} have cosine distance: {}".
 format(w1, w3, w1_w3_dist))
561
562 ### SOLUTION END

563 #%% md
564 ### SOLUTION BEGIN
565 A possible explanation for why this counter-
 intuitive result may have happened is that
 although synonyms share similar meanings, antonyms
 often co-occur in similar contexts and therefore
 share some semantic features that represented much
 more in the data than the synonyms.
566 In this case, light and heavy are both related to
 weight, and probably appear in the same context in
 the date more than light and illuminate.
567
568 ### SOLUTION END

569 #%% md
570 ### Question 2.4: Analogies with Word Vectors [written]
571 Word vectors have been shown to *sometimes*
 exhibit the ability to solve analogies.
572

```

573 As an example, for the analogy "man : grandfather :: woman : x" (read: man is to grandfather as woman is to x), what is x?

574

575 In the cell below, we show you how to use word vectors to find x using the `most\_similar` function from the \_\_[GenSim documentation]([https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.KeyedVectors.most\\_similar](https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.KeyedVectors.most_similar))\_\_. The function finds words that are most similar to the words in the `positive` list and most dissimilar from the words in the `negative` list (while omitting the input words, which are often the most similar; see [this paper](<https://www.aclweb.org/anthology/N18-2039.pdf>)). The answer to the analogy will have the highest cosine similarity (largest returned numerical value).

---

576 #%%

577 # Run this cell to answer the analogy -- man : grandfather :: woman : x

578 pprint.pprint(wv\_from\_bin.most\_similar(positive=['woman', 'grandfather'], negative=['man']))

---

579 #%% md

580 Let \$m\$, \$g\$, \$w\$, and \$x\$ denote the word vectors for `man`, `grandfather`, `woman`, and the answer, respectively. Using \*\*only\*\* vectors \$m\$, \$g\$, \$w\$, and the vector arithmetic operators \$+\$ and \$-\$ in your answer, to what expression are we maximizing \$x\$'s cosine similarity?

581

582 Hint: Recall that word vectors are simply multi-dimensional vectors that represent a word. It might help to draw out a 2D example using arbitrary locations of each vector. Where would `man` and `woman` lie in the coordinate plane relative to `grandfather` and the answer?

---

583 #%% md

584 ### SOLUTION BEGIN

585 Using vector arithmetic, we can express relationship between `man` & `grandfather` as:

```
586
587 g-m
588
589 This represents the vector that, when added to "man," results in "grandfather."
590
591 To find the analogous relationship for "woman," we will need add this difference vector to "woman":
592
593 w+(g-m)
594
595 Therefore, we can express X as:
596 x = w + (g - m)
597 ### SOLUTION END
598 #%% md
599 ### Question 2.5: Finding Analogies [code + written]
600 a. For the previous example, it's clear that "grandmother" completes the analogy. But give an intuitive explanation as to why the `most_similar` function gives us words like "granddaughter", "daughter", or "mother?"
601 #%% md
602 ### SOLUTION BEGIN
603 These following similar words probably occur in similar contexts which is why they are ranked highly. This is likely due to the fact that these words share similar semantic features, such as being related to familial relationships.
604 ### SOLUTION END
605 #%% md
606 b. Find an example of analogy that holds according to these vectors (i.e. the intended word is ranked top). In your solution please state the full analogy in the form x:y :: a:b. If you believe the analogy is complicated, explain why the analogy holds in one or two sentences.
607
608 **Note**: You may have to try many analogies to find one that works!
609 #%%
```

```

610 ### SOLUTION BEGIN
611
612 x = 'paris'
613 y = 'france'
614 a = 'berlin'
615 b = 'germany'
616 assert wv_from_bin.most_similar(positive=[a, y],
617 negative=[x])[0][0] == b
618 ### SOLUTION END
619 #%% md
620 ### SOLUTION BEGIN
621 The full analogy is: "paris : france :: berlin : germany".
622 This analogy holds because Paris is the capital city of France, just as Berlin is the capital city of Germany. The vectors for these words are likely to be close to each other in the vector space due to their shared semantic relationship as capital cities and countries.
623 ### SOLUTION END
624 #%% md
625 ### Question 2.6: Incorrect Analogy [code + written]
626 a. Below, we expect to see the intended analogy "hand : glove :: foot : **sock**", but we see an unexpected result instead. Give a potential reason as to why this particular analogy turned out the way it did?
627 #%%
628 pprint.pprint(wv_from_bin.most_similar(positive=['foot', 'glove'], negative=['hand']))
629 #%% md
630 ### SOLUTION BEGIN
631 A potential reason for this unexpected result that a 'foot' is a polysemous that often associated with a measurement of area (square feet).
632
633 The model has thus found words that describe areas instead of the expected association of "foot".
634

```

```

635 This could be due to the fact that the word "foot"
 " is more often used in contexts related to
 measurements, and the model has learned to
 associate it with other words that are commonly
 used in those contexts.
636 ### SOLUTION END
637 #%% md
638 b. Find another example of analogy that does *not*
 * hold according to these vectors. In your
 solution, state the intended analogy in the form x
 :y :: a:b, and state the **incorrect** value of b
 according to the word vectors (in the previous
 example, this would be **'45,000-square'**).
639 #%%
640 ### SOLUTION BEGIN
641
642 # correct analogy
643 x, y, a, b = "start", "stop", "execute", "cancel"
644 # incorrect analogy
645 pprint.pprint(wv_from_bin.most_similar(positive=[a
 , y], negative=[x]))
646
647 ### SOLUTION END
648 #%% md
649 ### SOLUTION BEGIN
650 The expected analogy is "start : stop :: execute
 : cancel", where "execute" is meant in the sense
 of "to launch" or "to carry out". However, the
 model instead returns words related to capital
 punishment, such as "kill" and "punish".
651 ### SOLUTION END
652 #%% md
653 ### Question 2.7: Guided Analysis of Bias in Word
 Vectors [written]
654
655 It's important to be cognizant of the biases (
 gender, race, sexual orientation etc.) implicit in
 our word embeddings. Bias can be dangerous
 because it can reinforce stereotypes through
 applications that employ these models.
656

```

```

657 Run the cell below, to examine (a) which terms are
 most similar to "woman" and "profession" and most
 dissimilar to "man", and (b) which terms are most
 similar to "man" and "profession" and most
 dissimilar to "woman". Point out the difference
 between the list of female-associated words and
 the list of male-associated words, and explain how
 it is reflecting gender bias.
658 #%%
659 # Run this cell
660 # Here `positive` indicates the list of words to
 be similar to and `negative` indicates the list of
 words to be
661 # most dissimilar from.
662
663 pprint.pprint(wv_from_bin.most_similar(positive=['
 man', 'profession'], negative=['woman']))
664 print()
665 pprint.pprint(wv_from_bin.most_similar(positive=['
 woman', 'profession'], negative=['man']))
666 #%% md
667 ### SOLUTION BEGIN
668 The difference between the two lists reveals a
 gender bias in how these word vectors represent
 the relationship between gender and professions.
 The women vector tends to contain more "soft" and
 specific words like: vocation and teaching, while
 the men vector contains more powerful and abstract
 words like: reputation, ethic, business.
669 ### SOLUTION END
670 #%% md
671 ### Question 2.8: Independent Analysis of Bias in
 Word Vectors [code + written]
672
673 Use the `most_similar` function to find another
 pair of analogies that demonstrates some bias is
 exhibited by the vectors. Please briefly explain
 the example of bias that you discover.
674 #%%
675 ### SOLUTION BEGIN
676

```

```
677 A = "caucasian"
678 B = "african"
679 word = "education"
680 pprint.pprint(wv_from_bin.most_similar(positive=[A
, word], negative=[B]))
681 print()
682 pprint.pprint(wv_from_bin.most_similar(positive=[B
, word], negative=[A]))
683
684 ### SOLUTION END
685 #%% md
686 ### SOLUTION BEGIN
687 The model associates "Caucasian + education" with
formal, academic concepts, and "African +
education" with development, aid, and public
health – reinforcing existing stereotypes about
educational access and systems based on race.
688 ### SOLUTION END
689 #%% md
690 ### Question 2.9: Thinking About Bias [written]
691
692 a. Give one explanation of how bias gets into the
word vectors. Briefly describe a real-world
example that demonstrates this source of bias.
693 #%% md
694 ### SOLUTION BEGIN
695 Bias gets into the word vectors primarily through
the data they are trained on. If the data contains
societal biases, those biases can be reflected in
the word vectors. For example, if the training
data frequently associates certain professions
with one gender, the word vectors will pick up on
these associations. The data might contain a
higher frequency of phrases like "female nurse"
and "male doctor," reflecting the historical
gender roles in these professions. As a result,
the trained word vectors might place "nurse"
closer to "woman" and "doctor" closer to "man" in
the embedding space, perpetuating and even
amplifying this gender bias.
696 ### SOLUTION END
```

```
697 #%% md
698 b. What is one method you can use to mitigate bias
 exhibited by word vectors? Briefly describe a
 real-world example that demonstrates this method.

699 #%% md
700
701 ### SOLUTION BEGIN
702 One method to mitigate this bias is to preprocess
 the data and modify biased occurrences or by
 augmenting the dataset with additional examples
 that counterbalance the bias.
703 One more method is to use debiasing algorithms
 that adjust the word vectors after training to
 reduce the bias.
704 A real world example of this is the "Word
 Embedding Association Test" (WEAT), which
 evaluates the bias in word embeddings by measuring
 the association between target words (e.g.,
 professions) and attribute words (e.g., gendered
 terms). By using this test, researchers can
 identify biased associations in word vectors and
 then apply debiasing techniques to reduce these
 biases.
705
706
707 ### SOLUTION END

708 #%% md
709 # Submission Instructions
710
711 1. Click the Save button at the top of the Jupyter
 Notebook.
712 2. Select Cell -> All Output -> Clear. This will
 clear all the outputs from all cells (but will
 keep the content of all cells).
713 2. Select Cell -> Run All. This will run all the
 cells in order, and will take several minutes.
714 3. Once you've rerun everything, select File ->
 Download as -> PDF via LaTeX (If you have trouble
 using "PDF via LaTeX", you can also save the
 webpage as pdf. Make sure all
```

- 714 your solutions especially the coding parts are displayed in the pdf</font>, it's okay if the provided codes get cut off because lines are not wrapped in code cells).
- 715 4. Look at the PDF file and make sure all your solutions are there, displayed correctly. The PDF is the only thing your graders will see!
- 716 5. Submit your PDF on \*\*MAMA\*\* system.