

# CSCI 4140: Natural Language Processing

## CSCI/DASC 6040: Computational Analysis of Natural Languages

Spring 2024

Homework 2 - Exploring word vectors

Due Sunday, February 4, at 11:59 PM

Do not redistribute without the instructor's written permission.

A lot of code is provided in this notebook, and we highly encourage you to read and understand it, as part of your learning.

**Assignment Notes:** Please make sure to save the notebook as you go along.

```
!pip install gensim
!pip install nltk
!pip install scikit-learn

Requirement already satisfied: gensim in /usr/local/lib/python3.10/dist-packages (4.3.2)
Requirement already satisfied: numpy>=1.18.5 in /usr/local/lib/python3.10/dist-packages (from gensim) (1.23.5)
Requirement already satisfied: scipy>=1.7.0 in /usr/local/lib/python3.10/dist-packages (from gensim) (1.11.4)
Requirement already satisfied: smart-open>=1.8.1 in /usr/local/lib/python3.10/dist-packages (from gensim) (6.4.0)
Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.8.1)
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.3.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2023.12.25)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.2.2)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.23.5)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.11.4)
Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.3.2)
Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.2.0)

# All Import Statements Defined Here
# Note: Do not add to this list.
# -------

import sys
assert sys.version_info[0]==3
assert sys.version_info[1] >= 5

from platform import python_version
assert int(python_version().split(".")[1]) >= 5, "Your Python version is " + python_version()

from gensim.models import KeyedVectors
from gensim.test.utils import datapath
import pprint
import matplotlib.pyplot as plt
plt.rcParams['figure.figsize'] = [10, 5]

import nltk
nltk.download('reuters') #to specify download location, optionally add the argument: download_dir='/specify/desired/path/'
from nltk.corpus import reuters

import numpy as np
import random
import scipy as sp
from sklearn.decomposition import TruncatedSVD
from sklearn.decomposition import PCA

START_TOKEN = '<START>'
END_TOKEN = '<END>'

np.random.seed(0)
random.seed(0)
# -------

[nltk_data] Downloading package reuters to /root/nltk_data...
```

## Word Vectors

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*.

**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](#) 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".

### Part 1: Count-Based Word Vectors (50 pts)

Most word vector models start from the following idea:

*You shall know a word by the company it keeps* ([Firth, J. R. 1957:11](#))

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 [here](#)).

#### Co-Occurrence

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.

**Example: Co-Occurrence with Fixed Window of n=1:**

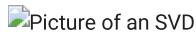
Document 1: "all that glitters is not gold"

Document 2: "all is well that ends well"

*	<START>	all	that	glitters	is	not	gold	well	ends	<END>
<START>	0	2	0	0	0	0	0	0	0	0
all	2	0	1	0	1	0	0	0	0	0
that	0	1	0	1	0	0	0	1	1	0
glitters	0	0	1	0	1	0	0	0	0	0
is	0	1	0	1	0	1	0	1	0	0
not	0	0	0	0	1	0	1	0	0	0
gold	0	0	0	0	0	1	0	0	0	1
well	0	0	1	0	1	0	0	0	1	1
ends	0	0	1	0	0	0	1	0	0	0
<END>	0	0	0	0	0	1	1	0	0	0

**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.

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  $S$  matrix, and our new, shorter length- $k$  word vectors in  $U_k$ .



This reduced-dimensionality co-occurrence representation preserves semantic relationships between words, e.g. *doctor* and *hospital* will be closer than *doctor* and *dog*.

**Notes:** If you can barely remember what an eigenvalue is, here's [a slow, friendly introduction to SVD](#). If you want to learn more thoroughly about PCA or SVD, feel free to check out lectures [7](#), [8](#), and [9](#) of CS168. These course notes provide a great high-level treatment of these general purpose algorithms. Though, for the purpose of this class, you only need to know how to extract the k-dimensional embeddings by utilizing pre-

programmed implementations of these algorithms from the numpy, scipy, or sklearn python packages. 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](#) – then there are reasonably scalable techniques to compute those iteratively.

## Plotting Co-Occurrence Word Embeddings

Here, we will be using the Reuters (business and financial news) corpus. If you haven't run the 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.

```
def read_corpus(category="gold"):
    """ Read files from the specified Reuter's category.
    Params:
        category (string): category name
    Return:
        list of lists, with words from each of the processed files
    """
    files = reuters.fileids(category)
    return [[START_TOKEN] + [w.lower() for w in list(reuters.words(f))] + [END_TOKEN] for f in files]
```

Let's have a look what these documents are like....

```
reuters_corpus = read_corpus()
pprint.pprint(reuters_corpus[:3], compact=True, width=100)

[[['<START>', 'western', 'mining', 'to', 'open', 'new', 'gold', 'mine', 'in', 'australia', 'western',
  'mining', 'corp', 'holdings', 'ltd', '&', 'lt', ';', 'wmng', '.', 's', '>', '(', 'wmc', ')',
  'said', 'it', 'will', 'establish', 'a', 'new', 'joint', 'venture', 'gold', 'mine', 'in', 'the',
  'northern', 'territory', 'at', 'a', 'cost', 'of', 'about', '21', 'mln', 'dlrs', '.', 'the',
  'mine', ',', 'to', 'be', 'known', 'as', 'the', 'goodall', 'project', ',', 'will', 'be', 'owned',
  '60', 'pct', 'by', 'wmc', 'and', '40', 'pct', 'by', 'a', 'local', 'w', '.', 'r', '.', 'grace',
  'and', 'co', '&', 'lt', ';', 'gra', '>', 'unit', '.', 'it', 'is', 'located', '30', 'kms', 'east',
  'of', 'the', 'adelaide', 'river', 'at', 'mt', '.', 'bundey', ',', 'wmc', 'said', 'in', 'a',
  'statement', 'it', 'said', 'the', 'open', '-', 'bit', 'mine', ',', 'with', 'a', 'conventional',
  'leach', 'treatment', 'plant', ',', 'is', 'expected', 'to', 'produce', 'about', '50', ',', '000',
  'ounces', 'of', 'gold', 'in', 'its', 'first', 'year', 'of', 'production', 'from', 'mid', '-',
  '1988', '.', 'annual', 'ore', 'capacity', 'will', 'be', 'about', '750', ',', '000', 'tonnes', '.',
  '<END>'],
 ['<START>', 'belgium', 'to', 'issue', 'gold', 'warrants', ',', 'sources', 'say', 'belgium',
  'plans', 'to', 'issue', 'swiss', 'franc', 'warrants', 'to', 'buy', 'gold', ',', 'with', 'credit',
  'suisse', 'as', 'lead', 'manager', ',', 'market', 'sources', 'said', '.', 'no', 'confirmation',
  'or', 'further', 'details', 'were', 'immediately', 'available', '.', '<END>'],
 ['<START>', 'belgium', 'launches', 'bonds', 'with', 'gold', 'warrants', 'the', 'kingdom', 'of',
  'belgium', 'is', 'launching', '100', 'mln', 'swiss', 'francs', 'of', 'seven', 'year', 'notes',
  'with', 'warrants', 'attached', 'to', 'buy', 'gold', ',', 'lead', 'mananger', 'credit', 'suisse',
  'said', '.', 'the', 'notes', 'have', 'a', '3', '-', '3', '/', '8', 'pct', 'coupon',
  'and', 'are', 'priced', 'at', 'par', '.', 'payment', 'is', 'due', 'april', '30', ',', '1987',
  'and', 'final', 'maturity', 'april', '30', ',', '1994', '.', 'each', '50', ',', '000', 'franc',
  'note', 'carries', '15', 'warrants', '.', 'two', 'warrants', 'are', 'required', 'to', 'allow',
  'the', 'holder', 'to', 'buy', '100', 'grammes', 'of', 'gold', 'at', 'a', 'price', 'of', '2', ',',
  '450', 'francs', ',', 'during', 'the', 'entire', 'life', 'of', 'the', 'bond', '.', 'the',
  'latest', 'gold', 'price', 'in', 'zurich', 'was', '2', ',', '045', '/', '2', ',', '070', 'francs',
  'per', '100', 'grammes', '.', '<END>']]
```

## Question 1.1: Implement `distinct_words` [code] (10 pts)

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 do it with Python list comprehensions. In particular, [this](#) may be useful to flatten a list of lists. If you're not familiar with Python list comprehensions in general, here's [more information](#).

Your returned `corpus_words` should be sorted. You can use python's `sorted` function for this.

You may find it useful to use [Python sets](#) to remove duplicate words.

```

def distinct_words(corpus):
    """ Determine a list of distinct words for the corpus.
    Params:
        corpus (list of list of strings): corpus of documents
    Return:
        corpus_words (list of strings): sorted list of distinct words across the corpus
        n_corpus_words (integer): number of distinct words across the corpus
    """
    # Flatten the list of lists and convert to lowercase
    flattened_corpus = [word for document in corpus for word in document]

    # Use set to get distinct words
    unique_words_set = set(flattened_corpus)

    # Sort the distinct words
    corpus_words = sorted(list(unique_words_set))

    # Number of distinct words
    n_corpus_words = len(corpus_words)

    return corpus_words, n_corpus_words

# -----
# Run this sanity check
# Note that this not an exhaustive check for correctness.
# -----


# Define toy corpus
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(" ")]
test_corpus_words, num_corpus_words = distinct_words(test_corpus)

# Correct answers
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_TOKEN])
ans_num_corpus_words = len(ans_test_corpus_words)

# Test correct number of words
assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of distinct words. Correct: {}. Yours: {}".format(ans_num_corpus_words, num_corpus_words)

# Test correct words
assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.\nCorrect: {}\nYours: {}".format(str(ans_test_corpus_words), str(test_corpus_words))

# Print Success
print ("-" * 80)
print("Passed All Tests!")
print ("-" * 80)

```

-----  
Passed All Tests!  
-----

### ▼ Question 1.2: Implement `compute_co_occurrence_matrix` [code] (15 pts)

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.

```
import numpy as np

def compute_co_occurrence_matrix(corpus, window_size=4):
    """ Compute co-occurrence matrix for the given corpus and window_size (default of 4).

    Note: Each word in a document should be at the center of a window. Words near edges will have a smaller
    number of co-occurring words.

    For example, if we take the document "<START> All that glitters is not gold <END>" with window size of 4,
    "All" will co-occur with "<START>", "that", "glitters", "is", and "not".

    Params:
        corpus (list of list of strings): corpus of documents
        window_size (int): size of context window
    Return:
        M (a symmetric numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus)):
            Co-occurrence matrix of word counts.
            The ordering of the words in the rows/columns should be the same as the ordering of the words given by the distinct_words fu
        word2ind (dict): dictionary that maps word to index (i.e. row/column number) for matrix M.
    """
    words, n_words = distinct_words(corpus)
    M = np.zeros((n_words, n_words))
    word2ind = {word: idx for idx, word in enumerate(words)}

    for doc in corpus:
        for i, target_word in enumerate(doc):
            target_index = word2ind[target_word]
            start = max(0, i - window_size)
            end = min(len(doc), i + window_size + 1)
            context_words = doc[start:i] + doc[i+1:end]

            for context_word in context_words:
                context_index = word2ind[context_word]
                M[target_index, context_index] += .5

    # Make the matrix symmetric
    M_symmetric = M + M.T - np.diag(M.diagonal())
    print(M_symmetric, word2ind)
    return M_symmetric, word2ind
```

```

# -----
# Run this sanity check
# Note that this is not an exhaustive check for correctness.
# -----


# Define toy corpus and get student's co-occurrence matrix
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(" ")]
M_test, word2ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)

# Correct M and word2ind
M_test_ans = np.array([
    [0., 0., 0., 0., 0., 0., 1., 0., 0., 1.],
    [0., 0., 1., 0., 0., 0., 0., 0., 0., 0.],
    [0., 1., 0., 0., 0., 0., 0., 1., 0., 0.],
    [0., 1., 0., 0., 0., 0., 0., 0., 0., 1.],
    [0., 0., 0., 0., 0., 0., 0., 0., 1., 0.],
    [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.],
    [1., 0., 0., 0., 0., 0., 0., 1., 0., 0.],
    [0., 0., 0., 0., 1., 1., 0., 0., 0., 0.],
    [0., 0., 1., 0., 1., 0., 0., 0., 0., 1.],
    [1., 0., 0., 1., 0., 0., 0., 1., 0., 0.]
])
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_TOKEN])
word2ind_ans = dict(zip(ans_test_corpus_words, range(len(ans_test_corpus_words))))


# Test correct word2ind
assert (word2ind_ans == word2ind_test), "Your word2ind is incorrect:\nCorrect: {}\nYours: {}".format(word2ind_ans, word2ind_test)

# Test correct M shape
assert (M_test.shape == M_test_ans.shape), "M matrix has incorrect shape.\nCorrect: {}\nYours: {}".format(M_test.shape, M_test_ans.shape)

# Test correct M values
for w1 in word2ind_ans.keys():
    idx1 = word2ind_ans[w1]
    for w2 in word2ind_ans.keys():
        idx2 = word2ind_ans[w2]
        student = M_test[idx1, idx2]
        correct = M_test_ans[idx1, idx2]
        if student != correct:
            print("Correct M:")
            print(M_test_ans)
            print("Your M: ")
            print(M_test)
            raise AssertionError("Incorrect count at index ({}, {})=({}, {}) in matrix M. Yours has {} but should have {}.".format(idx1, idx2, correct, student))

# Print Success
print("-" * 80)
print("Passed All Tests!")
print("-" * 80)

```

[[0. 0. 0. 0. 0. 1. 0. 0. 1.]  
[0. 0. 1. 1. 0. 0. 0. 0. 0.]  
[0. 1. 0. 0. 0. 0. 0. 1. 0.]  
[0. 1. 0. 0. 0. 0. 0. 0. 1.]  
[0. 0. 0. 0. 0. 0. 0. 1. 1.]  
[0. 0. 0. 0. 0. 0. 1. 1. 0.]  
[1. 0. 0. 0. 0. 0. 1. 0. 0.]  
[0. 0. 0. 0. 1. 1. 0. 0. 0.]  
[0. 0. 1. 0. 1. 1. 0. 0. 1.]  
[1. 0. 0. 1. 1. 0. 0. 1. 0.]] {'<END>': 0, '<START>': 1, 'All': 2, "All's": 3, 'ends': 4, 'glitters': 5, 'gold': 6, "isn't": 7, 'tha': 8}

---

Passed All Tests!

### ▼ Question 1.3: Implement reduce\_to\_k\_dim [code] (5 pts)

Construct a method that performs dimensionality reduction on the matrix to produce k-dimensional embeddings. Use SVD to take the top k components and produce a new matrix of k-dimensional embeddings.

**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](#).

```

def reduce_to_k_dim(M, k=2):
    """ Reduce a co-occurrence count matrix of dimensionality (num_corpus_words, num_corpus_words)
        to a matrix of dimensionality (num_corpus_words, k) using the following SVD function from Scikit-Learn:
        - http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html

    Params:
        M (numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus)): co-occurrence matrix of
        k (int): embedding size of each word after dimension reduction
    Returns:
        M_reduced (numpy matrix of shape (number of corpus words, k)): matrix of k-dimensioal word embeddings.
        In terms of the SVD from math class, this actually returns U * S
    """
    n_iters = 10 # Use this parameter in your call to `TruncatedSVD`
    svd = TruncatedSVD(n_components=k, n_iter=n_iters)
    M_reduced = svd.fit_transform(M)

    print("Done.")
    return M_reduced

# -----
# Run this sanity check
# Note that this is not an exhaustive check for correctness
# In fact we only check that your M_reduced has the right dimensions.
# -----

# Define toy corpus and run student code
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(" ")]
M_test, word2ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)
M_test_reduced = reduce_to_k_dim(M_test, k=2)

# Test proper dimensions
assert (M_test_reduced.shape[0] == 10), "M_reduced has {} rows; should have {}".format(M_test_reduced.shape[0], 10)
assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have {}".format(M_test_reduced.shape[1], 2)

# Print Success
print ("-" * 80)
print("Passed All Tests!")
print ("-" * 80)

[[0. 0. 0. 0. 0. 0. 1. 0. 0. 1.]
 [0. 0. 1. 1. 0. 0. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0. 1. 0.]
 [0. 1. 0. 0. 0. 0. 0. 0. 0. 1.]
 [0. 0. 0. 0. 0. 0. 0. 1. 1.]
 [0. 0. 0. 0. 0. 0. 1. 1. 0.]
 [0. 0. 0. 0. 0. 0. 1. 0. 0.]
 [1. 0. 0. 0. 0. 0. 0. 1. 0. 0.]
 [0. 0. 0. 0. 0. 1. 1. 0. 0. 0.]
 [0. 0. 1. 0. 1. 1. 0. 0. 0. 1.]
 [1. 0. 0. 1. 1. 0. 0. 1. 0. 0.]] {'<END>': 0, '<START>': 1, 'All': 2, "All's": 3, 'ends': 4, 'glitters': 5, 'gold': 6, "isn't": 7, 'tha
Done.

-----
Passed All Tests!
-----
```

#### ✓ Question 1.4: Implement plot\_embeddings [code] (5 pts)

Here you will write a function to plot a set of 2D vectors in 2D space. For graphs, we will use Matplotlib (`plt`).

For this example, you may find it useful to adapt [this code](#). In the future, a good way to make a plot is to look at [the Matplotlib gallery](#), find a plot that looks somewhat like what you want, and adapt the code they give.

```
def plot_embeddings(M_reduced, word2ind, words):
    """ Plot in a scatterplot the embeddings of the words specified in the list "words".
    NOTE: do not plot all the words listed in M_reduced / word2ind.
    Include a label next to each point.

    Params:
        M_reduced (numpy matrix of shape (number of unique words in the corpus , 2)): matrix of 2-dimensioal word embeddings
        word2ind (dict): dictionary that maps word to indices for matrix M
        words (list of strings): words whose embeddings we want to visualize
    """

    # Get the indices of the specified words
    indices = [word2ind[word] for word in words]

    # Extract the coordinates of the specified words
    x_coords = M_reduced[indices, 0]
    y_coords = M_reduced[indices, 1]

    # Plot the points
    plt.figure(figsize=(10, 8))
    plt.scatter(x_coords, y_coords, color='blue')

    # Annotate each point with its word label
    for label, x, y in zip(words, x_coords, y_coords):
        plt.annotate(label, xy=(x, y), xytext=(5, 2), textcoords='offset points', ha='right', va='bottom')

    plt.title('Word Embeddings Visualization')
    plt.xlabel('Dimension 1')
    plt.ylabel('Dimension 2')
    plt.show()

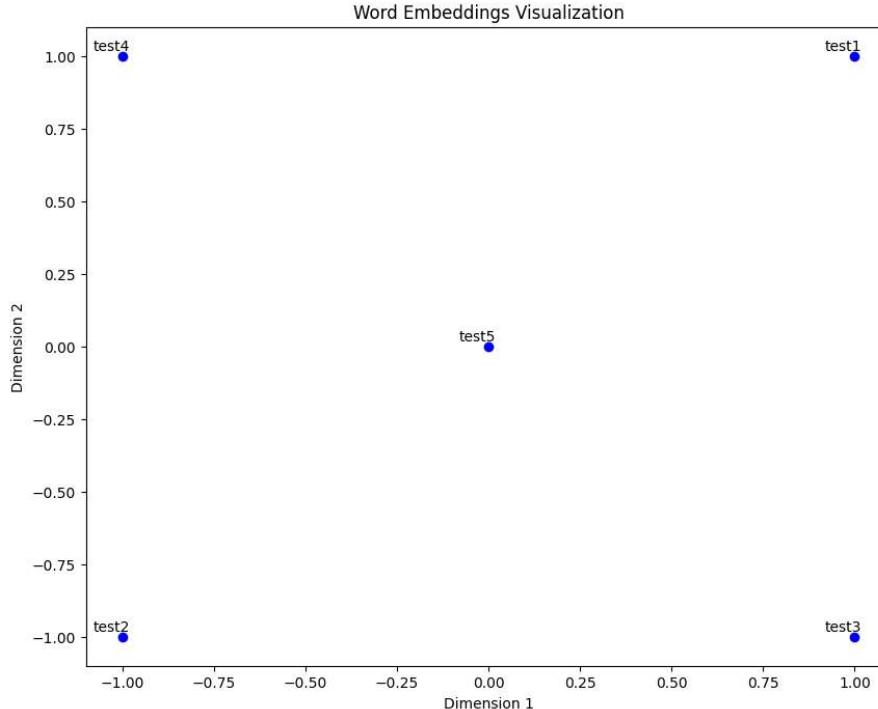
# -----
# Run this sanity check
# Note that this is not an exhaustive check for correctness.
# The plot produced should look like the "test solution plot" depicted below.
# -----"

print ("-" * 80)
print ("Outputted Plot:")

M_reduced_plot_test = np.array([[1, 1], [-1, -1], [1, -1], [-1, 1], [0, 0]])
word2ind_plot_test = {'test1': 0, 'test2': 1, 'test3': 2, 'test4': 3, 'test5': 4}
words = ['test1', 'test2', 'test3', 'test4', 'test5']
plot_embeddings(M_reduced_plot_test, word2ind_plot_test, words)

print ("-" * 80)
```

Outputted Plot:



#### ✓ Question 1.5: Co-Occurrence Plot Analysis [written] (15 pts)

Now we will put together all the parts you have written! We will compute the co-occurrence matrix with fixed window of 4 (the default window size), over the Reuters "gold" corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns  $U^*S$ , so we need to normalize the returned vectors, so that all the vectors will appear around the unit circle (therefore closeness is directional closeness). **Note:** The line of code below that does the normalizing uses the NumPy concept of *broadcasting*. If you don't know about broadcasting, check out [Computation on Arrays: Broadcasting by Jake VanderPlas](#).

Run the below cell to produce the plot. It'll probably take a few seconds to run.

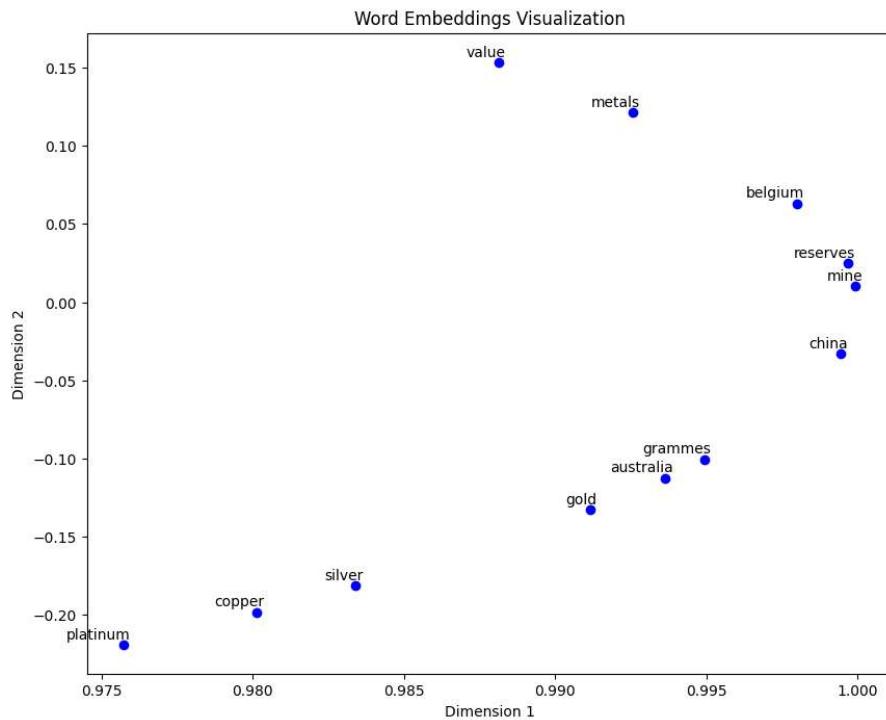
```
# -----
# Run This Cell to Produce Your Plot
# -----
reuters_corpus = read_corpus()
M_co_occurrence, word2ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

# Rescale (normalize) the rows to make them each of unit-length
M_lengths = np.linalg.norm(M_reduced_co_occurrence, axis=1)
M_normalized = M_reduced_co_occurrence / M_lengths[:, np.newaxis] # broadcasting

words = ['value', 'gold', 'platinum', 'reserves', 'silver', 'metals', 'copper', 'belgium', 'australia', 'china', 'grammes', "mine"]

plot_embeddings(M_normalized, word2ind_co_occurrence, words)
```

```
[[11.  0.  3. ... 0.  3.  0.]
 [ 0.  0.  6. ... 0.  0.  0.]
 [ 3.  6.  2. ... 0.  0.  0.]
 ...
 [ 0.  0.  0. ... 0.  0.  0.]
 [ 3.  0.  0. ... 0.  0.  0.]
 [ 0.  0.  0. ... 0.  0.  0.]] {'"': 0, '&': 1, '"": 2, '(': 3, ')': 4, ',': 5, '.': 6}
Done.
```



Verify that your figure matches "question\_1.5.png" in the assignment zip. If not, use that figure to answer the next two questions.



- a. Find at least two groups of words that cluster together in 2-dimensional embedding space. Give an explanation for each cluster you observe.

**Answer:**

Copper, platinum, and australias, belgium are clustered together. This is because these terms are related and appear together more often than say China and value.

- b. What doesn't cluster together that you might think should have? Describe at least two examples.

**Answer:**

I believe China should have clustered with belgium and australias. These are all related terms being places and country's. Also I believe gold and metal should have been clustered closer to platinum and copper and silver being they are all ore and metals.

✓ Part 2: Prediction-Based Word Vectors (60 pts)

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.

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.

```
def load_embedding_model():
    """ Load GloVe Vectors
    Return:
        wv_from_bin: All 400000 embeddings, each length 200
    """
    import gensim.downloader as api
    wv_from_bin = api.load("glove-wiki-gigaword-200")
    print("Loaded vocab size %i" % len(list(wv_from_bin.index_to_key)))
    return wv_from_bin

# -----
# Run Cell to Load Word Vectors
# Note: This will take a couple minutes
# -----
wv_from_bin = load_embedding_model()

[=====] 100.0% 252.1/252.1MB downloaded
Loaded vocab size 400000
```

**Note:** If you are receiving a "reset by peer" error, rerun the cell to restart the download. If you run into an "attribute" error, you may need to update to the most recent version of gensim and numpy. You can upgrade them inline by uncommenting and running the below cell:

```
#!pip install gensim --upgrade
#!pip install numpy --upgrade
```

## ✓ Reducing dimensionality of Word Embeddings

Let's directly compare the GloVe embeddings to those of the co-occurrence matrix. In order to avoid running out of memory, we will work with a sample of 10000 GloVe vectors instead. Run the following cells to:

1. Put 10000 Glove vectors into a matrix M
2. Run `reduce_to_k_dim` (your Truncated SVD function) to reduce the vectors from 200-dimensional to 2-dimensional.

```

def get_matrix_of_vectors(wv_from_bin, required_words):
    """ Put the GloVe vectors into a matrix M.
    Param:
        wv_from_bin: KeyedVectors object; the 400000 GloVe vectors loaded from file
    Return:
        M: numpy matrix shape (num words, 200) containing the vectors
        word2ind: dictionary mapping each word to its row number in M
    """
    import random
    words = list(wv_from_bin.index_to_key)
    print("Shuffling words ...")
    random.seed(225)
    random.shuffle(words)
    words = words[:10000]
    print("Putting %i words into word2ind and matrix M..." % len(words))
    word2ind = {}
    M = []
    curInd = 0
    for w in words:
        try:
            M.append(wv_from_bin.get_vector(w))
            word2ind[w] = curInd
            curInd += 1
        except KeyError:
            continue
    for w in required_words:
        if w in words:
            continue
        try:
            M.append(wv_from_bin.get_vector(w))
            word2ind[w] = curInd
            curInd += 1
        except KeyError:
            continue
    M = np.stack(M)
    print("Done.")
    return M, word2ind

# -----
# Run Cell to Reduce 200-Dimensional Word Embeddings to k Dimensions
# Note: This should be quick to run
# -----
M, word2ind = get_matrix_of_vectors(wv_from_bin, words)
M_reduced = reduce_to_k_dim(M, k=2)

# Rescale (normalize) the rows to make them each of unit-length
M_lengths = np.linalg.norm(M_reduced, axis=1)
M_reduced_normalized = M_reduced / M_lengths[:, np.newaxis] # broadcasting

Shuffling words ...
Putting 10000 words into word2ind and matrix M...
Done.
Done.

```

**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, try running this notebook on [Google Colab](#).

#### ✓ Question 2.1: GloVe Plot Analysis [written] (12 pts)

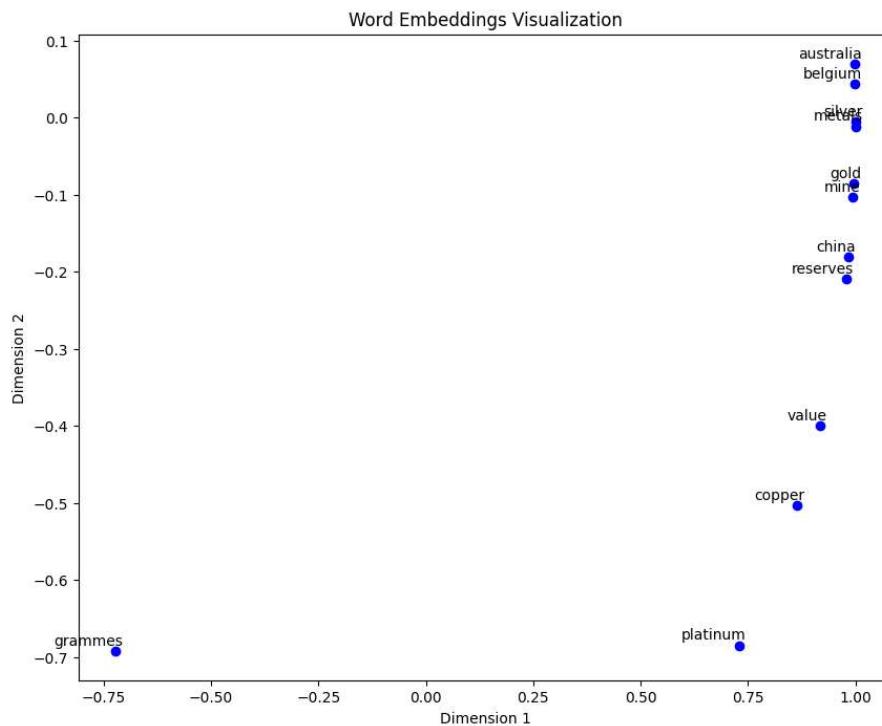
Run the cell below to plot the 2D GloVe embeddings for ['value', 'gold', 'platinum', 'reserves', 'silver', 'metals', 'copper', 'belgium', 'australia', 'china', 'grammes', 'mine'].

```

words = ['value', 'gold', 'platinum', 'reserves', 'silver', 'metals', 'copper',
'belgium', 'australia', 'china', 'grammes', "mine"]

plot_embeddings(M_reduced_normalized, word2ind, words)

```



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?

**Answer:**

the Dimensions are different. -.7 to .1 and -.75 to 1. Also the words are not in the same places and have different clusters. Some clusters are the same like belgium and australia

b. What is a possible cause for the difference?

**Answer:**

The sample/sample-size of words has changed

Double-click (or enter) to edit

## Cosine Similarity

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.

We can think of n-dimensional vectors as points in n-dimensional space. If we take this perspective L1 and L2 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:



Instead of computing the actual angle, we can leave the similarity in terms of  $similarity = \cos(\Theta)$ . Formally the [Cosine Similarity](#)  $s$  between two vectors  $p$  and  $q$  is defined as:

$$s = \frac{p \cdot q}{\|p\| \|q\|}, \text{ where } s \in [-1, 1]$$

### ✓ Question 2.2: Words with Multiple Meanings [code + written] (6 pts)

Polysemes and homonyms are words that have more than one meaning. 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.

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)?

**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](#).

```
# Let's find a word with multiple meanings
polysemous_word = "run"

# Get the top-10 most similar words
similar_words = wv_from_bin.most_similar(polysemous_word)

# Extract the meanings from the top-10
meanings = {
    'meaning_1': ['flying_mammal', 'winged_mammal', 'nocturnal_animal'],
    'meaning_2': ['sports_equipment', 'baseball', 'cricket', 'racket']
}

# Check if both meanings are present in the top-10
found_meanings = {key: any(word in [similar[0] for similar in similar_words] for word in meanings[key]) for key in meanings}
print(similar_words)

polysemous_word, found_meanings
```

[('running', 0.7378345727920532), ('runs', 0.7364052534103394), ('ran', 0.696038544178009), ('went', 0.6395289897918701), ('start', 0.63
('run', {'meaning\_1': False, 'meaning\_2': False})

#### Answer:

Run has different meanings. For example: running, runs, ran, these meanings all have to do with movement and actively using your body to RUN went, start, go, going, first are all a where meaning. where did you go? i went... im going... I go... I start...

The first word I tried was bat. Believing it would give me a bat like baseball and bat like the flying creature , but it only gave my synonyms for bat as in baseball. Maybe its because baseball is a lot more popular than bats as a creature are, because of this there are more examples of bats being associated with baseball.

### ✓ Question 2.3: Synonyms & Antonyms [code + written] (8 pts)

When considering Cosine Similarity, it's often more convenient to think of Cosine Distance, which is simply  $1 - \text{Cosine Similarity}$ .

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

As an example,  $w_1 = \text{"happy"}$  is closer to  $w_3 = \text{"sad"}$  than to  $w_2 = \text{"cheerful"}$ . 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.

You should use the the `wv_from_bin.distance(w1, w2)` function here in order to compute the cosine distance between two words. Please see the [GenSim documentation](#) for further assistance.

```
### SOLUTION BEGIN
```

```
w1 = "big"
w2 = "massive"
w3 = "small"
w1_w2_dist = wv_from_bin.distance(w1, w2)
w1_w3_dist = wv_from_bin.distance(w1, w3)

print("Synonyms {}, {} have cosine distance: {}".format(w1, w2, w1_w2_dist))
print("Antonyms {}, {} have cosine distance: {}".format(w1, w3, w1_w3_dist))
```

```
### SOLUTION END
```

```
Synonyms big, massive have cosine distance: 0.48080557584762573
Antonyms big, small have cosine distance: 0.3511464595794678
```

### Answer:

w1= big w2= massive w3= small. massive is a synonym to big but when you describe something as massive you mean extremely big. A closer synonym would be large. The exact antonym of big is small, this is why the cosine distance is smaller.

### ✓ Question 2.4: Analogies with Word Vectors [written] (6 pts)

Word vectors have been shown to sometimes exhibit the ability to solve analogies.

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

In the cell below, we show you how to use word vectors to find x using the `most_similar` function from the [GenSim documentation](#). 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](#)). The answer to the analogy will have the highest cosine similarity (largest returned numerical value).

```
# Run this cell to answer the analogy -- man : grandfather :: woman : x
pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'grandfather'], negative=['man']))

[('grandmother', 0.7608445286750793),
 ('granddaughter', 0.7200808525085449),
 ('daughter', 0.7168302536010742),
 ('mother', 0.7151536345481873),
 ('niece', 0.7005682587623596),
 ('father', 0.6659887433052063),
 ('aunt', 0.6623498794403076),
 ('grandson', 0.6618767976760864),
 ('grandparents', 0.644661009311676),
 ('wife', 0.6445354223251343)]
```

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?

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?

### Answer:

Let `man` = (1,1) `Woman` = (4,3) `Grandfather` = (6,6)

then using  $x = g + (w-m)$

$x = (9,8)$

### ✓ Question 2.5: Finding Analogies [code + written] (6 pts)

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"?

### Answer:

the first analogy gives a clue to what the answer for the second analogy is. The first analogy associates female to male, and young to old. Using this association the second analogy can predict male to female, and young to old, Man to grandmother. Also because grandmother is similar to woman and grandfather and just not similar enough to man

from the top of my head, the function sees man to grandfather and associates the words with male family members. Then it sees woman and using the previous assumptions female family members pops out.

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.

**Note:** You may have to try many analogies to find one that works!

```
### SOLUTION BEGIN
```

```
x, y, a, b = 'finger', 'toe', 'fingernail', 'toenail'
pprint.pprint(wv_from_bin.most_similar(positive=[a, y], negative=[x]))
assert wv_from_bin.most_similar(positive=[a,y], negative=[x])[0][0] == b
```

```
### SOLUTION END
```

```
[('toenail', 0.5517504215240479),
 ('chador', 0.4644639790058136),
 ('stiletto', 0.46429678797721863),
 ('stilettos', 0.457084596157074),
 ('tights', 0.45212215185165405),
 ('fishnet', 0.4364977478981018),
 ('pantyhose', 0.43322479724884033),
 ('scrapings', 0.42653197050094604),
 ('fishnets', 0.4066351354122162),
 ('strappy', 0.40389060974121894)]
```

**Answer:**

Finger:Fingernail::toe:toenail the analogy holds because it is specific and there is not a lot of other options that are similar.

#### ✓ Question 2.6: Incorrect Analogy [code + written] (6 pts)

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?

```
pprint.pprint(wv_from_bin.most_similar(positive=['foot', 'glove'], negative=['hand']))

[('45,000-square', 0.4922032654285431),
 ('15,000-square', 0.4649604558944702),
 ('10,000-square', 0.4544755816459656),
 ('6,000-square', 0.44975775480270386),
 ('3,500-square', 0.444133460521698),
 ('700-square', 0.44257497787475586),
 ('50,000-square', 0.4356396794319153),
 ('3,000-square', 0.43486514687538147),
 ('30,000-square', 0.4330596923828125),
 ('footed', 0.43236875534057617)]
```

**Answer:**

I believe it could be because foot has multiple meanings. Foot as in your feet or foot as in a unit of measure. Also, the correct association might be foot to shoe

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**').

```
### SOLUTION BEGIN
```