

exploring_word_vectors

December 10, 2019

1 Lab4: Word Vectors

2 Yucheng Jin (yucheng9)

```
[1]: # All Import Statements Defined Here
# Note: Do not add to this list.
# All the dependencies you need, can be installed by running .
# Make sure to have scikit-learn version 0.21.3, otherwise you might run into
# →import problems.
# I found it best to create a new conda environment with scikit-learn 0.21.3
# →installed using the command:
# conda create -n lab4_env python=3 scikit-learn==0.21.3
# -----

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

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')
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
[nltk_data] /Users/yuchengjin/nltk_data...
[nltk_data] Package reuters is already up-to-date!
```

2.1 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 *word2vec*.

Assignment Notes: Please make sure to save the notebook as you go along. Submission Instructions are located at the bottom of the notebook.

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*”.

2.2 Part 1: Count-Based Word Vectors (50 points)

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](#)).

2.2.1 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.

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

Document 1: “all that glitters is not gold”

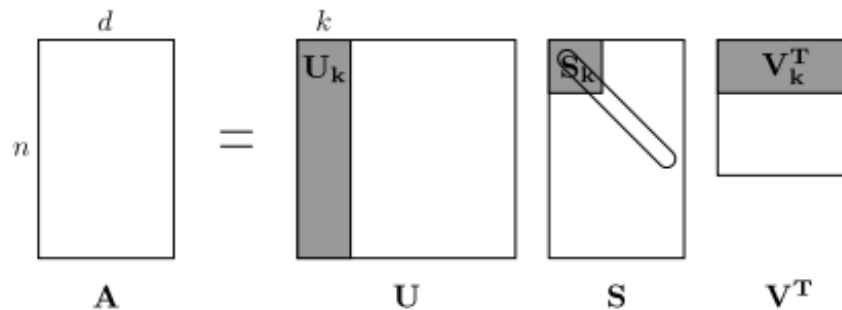
Document 2: “all is well that ends well”

*	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

*	START	all	that	glitters	is	not	gold	well	ends	END
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	0	1	0	0
END	0	0	0	0	0	0	1	1	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 .



Picture of an SVD

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](#). 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.

2.2.2 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 "crude" (i.e. news articles about oil, gas, 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.

```
[2]: def read_corpus(category="crude"):
    """ 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...

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

```
[[ '<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand',
  'downwards', 'the',
  'ministry', 'of', 'international', 'trade', 'and', 'industry', '(', 'miti',
  ')', 'will', 'revise',
  'its', 'long', '-', 'term', 'energy', 'supply', '/', 'demand', 'outlook',
  'by', 'august', 'to',
  'meet', 'a', 'forecast', 'downtrend', 'in', 'japanese', 'energy', 'demand',
  ', ', 'ministry',
  'officials', 'said', '.', 'miti', 'is', 'expected', 'to', 'lower', 'the',
  'projection', 'for',
  'primary', 'energy', 'supplies', 'in', 'the', 'year', '2000', 'to', '550',
  'mln', 'kilolitres',
  '(', 'kl', ')', 'from', '600', 'mln', ', ', 'they', 'said', '.', 'the',
  'decision', 'follows',
  'the', 'emergence', 'of', 'structural', 'changes', 'in', 'japanese',
  'industry', 'following',
  'the', 'rise', 'in', 'the', 'value', 'of', 'the', 'yen', 'and', 'a',
  'decline', 'in', 'domestic',
  'electric', 'power', 'demand', '.', 'miti', 'is', 'planning', 'to', 'work',
  'out', 'a', 'revised',
  'energy', 'supply', '/', 'demand', 'outlook', 'through', 'deliberations',
  'of', 'committee',
  'meetings', 'of', 'the', 'agency', 'of', 'natural', 'resources', 'and',
```

'energy', ',', 'the',
 'officials', 'said', '.', 'they', 'said', 'miti', 'will', 'also', 'review',
 'the', 'breakdown',
 'of', 'energy', 'supply', 'sources', ',', 'including', 'oil', ',', 'nuclear',
 ',', 'coal', 'and',
 'natural', 'gas', '.', 'nuclear', 'energy', 'provided', 'the', 'bulk', 'of',
 'japan', '"', 's',
 'electric', 'power', 'in', 'the', 'fiscal', 'year', 'ended', 'march', '31',
 ',', 'supplying',
 'an', 'estimated', '27', 'pct', 'on', 'a', 'kilowatt', '/', 'hour', 'basis',
 ',', 'followed',
 'by', 'oil', '(', '23', 'pct', ')', 'and', 'liquefied', 'natural', 'gas', '(',
 '21', 'pct', ')',
 'they', 'noted', '.', '<END>'],
 ['<START>', 'energy', '/', 'u', '.', 's', '.', 'petrochemical', 'industry',
 'cheap', 'oil',
 'feedstocks', ',', 'the', 'weakened', 'u', '.', 's', '.', 'dollar', 'and',
 'a', 'plant',
 'utilization', 'rate', 'approaching', '90', 'pct', 'will', 'propel', 'the',
 'streamlined', 'u',
 '.', 's', '.', 'petrochemical', 'industry', 'to', 'record', 'profits', 'this',
 'year', ',',
 'with', 'growth', 'expected', 'through', 'at', 'least', '1990', ',', 'major',
 'company',
 'executives', 'predicted', '.', 'this', 'bullish', 'outlook', 'for',
 'chemical', 'manufacturing',
 'and', 'an', 'industrywide', 'move', 'to', 'shed', 'unrelated', 'businesses',
 'has', 'prompted',
 'gaf', 'corp', '&', 'lt', ';', 'gaf', '>', 'privately', '-', 'held', 'cain',
 'chemical', 'inc',
 ',', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions',
 'of', 'petrochemical',
 'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&',
 'lt', ';', 'ash',
 '>', 'the', 'kentucky', '-', 'based', 'oil', 'refiner', 'and', 'marketer',
 ',', 'are', 'also',
 'shopping', 'for', 'money', '-', 'making', 'petrochemical', 'businesses',
 'to', 'buy', '.', '"',
 'i', 'see', 'us', 'poised', 'at', 'the', 'threshold', 'of', 'a', 'golden',
 'period', ',', '"', 'said',
 'paul', 'oreffice', ',', 'chairman', 'of', 'giant', 'dow', 'chemical', 'co',
 '&', 'lt', ';',
 'dow', '>', 'adding', ',', '"', 'there', '"', 's', 'no', 'major', 'plant',
 'capacity', 'being',
 'added', 'around', 'the', 'world', 'now', '.', 'the', 'whole', 'game', 'is',
 'bringing', 'out',
 'new', 'products', 'and', 'improving', 'the', 'old', 'ones', '.', 'analysts',
 'say', 'the',

'chemical', 'industry', '"', 's', 'biggest', 'customers', ',', 'automobile',
 'manufacturers',
 'and', 'home', 'builders', 'that', 'use', 'a', 'lot', 'of', 'paints', 'and',
 'plastics', ',',
 'are', 'expected', 'to', 'buy', 'quantities', 'this', 'year', '.', 'u', '.',
 's', '.',
 'petrochemical', 'plants', 'are', 'currently', 'operating', 'at', 'about',
 '90', 'pct',
 'capacity', ',', 'reflecting', 'tighter', 'supply', 'that', 'could', 'hike',
 'product', 'prices',
 'by', '30', 'to', '40', 'pct', 'this', 'year', ',', 'said', 'john', 'dosher',
 ',', 'managing',
 'director', 'of', 'pace', 'consultants', 'inc', 'of', 'houston', '.',
 'demand', 'for', 'some',
 'products', 'such', 'as', 'styrene', 'could', 'push', 'profit', 'margins',
 'up', 'by', 'as',
 'much', 'as', '300', 'pct', ',', 'he', 'said', '.', 'oreffice', ',',
 'speaking', 'at', 'a',
 'meeting', 'of', 'chemical', 'engineers', 'in', 'houston', ',', 'said', 'dow',
 'would', 'easily',
 'top', 'the', '741', 'mln', 'dlrs', 'it', 'earned', 'last', 'year', 'and',
 'predicted', 'it',
 'would', 'have', 'the', 'best', 'year', 'in', 'its', 'history', '.', 'in',
 '1985', ',', 'when',
 'oil', 'prices', 'were', 'still', 'above', '25', 'dlrs', 'a', 'barrel', 'and',
 'chemical',
 'exports', 'were', 'adversely', 'affected', 'by', 'the', 'strong', 'u', '.',
 's', '.', 'dollar',
 ',', 'dow', 'had', 'profits', 'of', '58', 'mln', 'dlrs', '.', '"', 'i',
 'believe', 'the',
 'entire', 'chemical', 'industry', 'is', 'headed', 'for', 'a', 'record',
 'year', 'or', 'close',
 'to', 'it', ',', '"', 'oreffice', 'said', '.', 'gaf', 'chairman', 'samuel',
 'heyman', 'estimated',
 'that', 'the', 'u', '.', 's', '.', 'chemical', 'industry', 'would', 'report',
 'a', '20', 'pct',
 'gain', 'in', 'profits', 'during', '1987', '.', 'last', 'year', ',', 'the',
 'domestic',
 'industry', 'earned', 'a', 'total', 'of', '13', 'billion', 'dlrs', ',', 'a',
 '54', 'pct', 'leap',
 'from', '1985', '.', 'the', 'turn', 'in', 'the', 'fortunes', 'of', 'the',
 'once', '-', 'sickly',
 'chemical', 'industry', 'has', 'been', 'brought', 'about', 'by', 'a',
 'combination', 'of', 'luck',
 'and', 'planning', ',', 'said', 'pace', '"', 's', 'john', 'dosher', '.',
 'dosher', 'said', 'last',
 'year', '"', 's', 'fall', 'in', 'oil', 'prices', 'made', 'feedstocks',
 'dramatically', 'cheaper',

'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was',
 'weakening', 'against',
 'foreign', 'currencies', '.', 'that', 'helped', 'boost', 'u', '.', 's', '.',
 'chemical',
 'exports', '.', 'also', 'helping', 'to', 'bring', 'supply', 'and', 'demand',
 'into', 'balance',
 'has', 'been', 'the', 'gradual', 'market', 'absorption', 'of', 'the', 'extra',
 'chemical',
 'manufacturing', 'capacity', 'created', 'by', 'middle', 'eastern', 'oil',
 'producers', 'in',
 'the', 'early', '1980s', '.', 'finally', ',', 'virtually', 'all', 'major',
 'u', '.', 's', '.',
 'chemical', 'manufacturers', 'have', 'embarked', 'on', 'an', 'extensive',
 'corporate',
 'restructuring', 'program', 'to', 'mothball', 'inefficient', 'plants', ',',
 'trim', 'the',
 'payroll', 'and', 'eliminate', 'unrelated', 'businesses', '.', 'the',
 'restructuring', 'touched',
 'off', 'a', 'flurry', 'of', 'friendly', 'and', 'hostile', 'takeover',
 'attempts', '.', 'gaf', ',',
 'which', 'made', 'an', 'unsuccessful', 'attempt', 'in', '1985', 'to',
 'acquire', 'union',
 'carbide', 'corp', '&', 'lt', ';', 'uk', '>', 'recently', 'offered', 'three',
 'billion', 'dlrs',
 'for', 'borg', 'warner', 'corp', '&', 'lt', ';', 'bor', '>', 'a', 'chicago',
 'manufacturer',
 'of', 'plastics', 'and', 'chemicals', '.', 'another', 'industry',
 'powerhouse', ',', 'w', '.',
 'r', '.', 'grace', '&', 'lt', ';', 'gra', '>', 'has', 'divested', 'its',
 'retailing', ',',
 'restaurant', 'and', 'fertilizer', 'businesses', 'to', 'raise', 'cash', 'for',
 'chemical',
 'acquisitions', '.', 'but', 'some', 'experts', 'worry', 'that', 'the',
 'chemical', 'industry',
 'may', 'be', 'headed', 'for', 'trouble', 'if', 'companies', 'continue',
 'turning', 'their',
 'back', 'on', 'the', 'manufacturing', 'of', 'staple', 'petrochemical',
 'commodities', ',', 'such',
 'as', 'ethylene', ',', 'in', 'favor', 'of', 'more', 'profitable', 'specialty',
 'chemicals',
 'that', 'are', 'custom', '-', 'designed', 'for', 'a', 'small', 'group', 'of',
 'buyers', '.', '""',
 'companies', 'like', 'dupont', '&', 'lt', ';', 'dd', '>', 'and', 'monsanto',
 'co', '&', 'lt', ';',
 'mtc', '>', 'spent', 'the', 'past', 'two', 'or', 'three', 'years', 'trying',
 'to', 'get', 'out',
 'of', 'the', 'commodity', 'chemical', 'business', 'in', 'reaction', 'to',
 'how', 'badly', 'the',

'market', 'had', 'deteriorated', ',,"', 'dosher', 'said', '.', '"', 'but', 'i',
 'think', 'they',
 'will', 'eventually', 'kill', 'the', 'margins', 'on', 'the', 'profitable',
 'chemicals', 'in',
 'the', 'niche', 'market', '."', 'some', 'top', 'chemical', 'executives',
 'share', 'the',
 'concern', '.', '"', 'the', 'challenge', 'for', 'our', 'industry', 'is', 'to',
 'keep', 'from',
 'getting', 'carried', 'away', 'and', 'repeating', 'past', 'mistakes', ',,"',
 'gaf', '"', 's',
 'heyman', 'cautioned', '.', '"', 'the', 'shift', 'from', 'commodity',
 'chemicals', 'may', 'be',
 'ill', '-', 'advised', '.', 'specialty', 'businesses', 'do', 'not', 'stay',
 'special', 'long',
 '."', 'houston', '-', 'based', 'cain', 'chemical', ',,', 'created', 'this',
 'month', 'by', 'the',
 'sterling', 'investment', 'banking', 'group', ',,', 'believes', 'it', 'can',
 'generate', '700',
 'mln', 'dlrs', 'in', 'annual', 'sales', 'by', 'bucking', 'the', 'industry',
 'trend', '.',
 'chairman', 'gordon', 'cain', ',,', 'who', 'previously', 'led', 'a',
 'leveraged', 'buyout', 'of',
 'dupont', '"', 's', 'conoco', 'inc', '"', 's', 'chemical', 'business', ',,',
 'has', 'spent', '1',
 '.', '1', 'billion', 'dlrs', 'since', 'january', 'to', 'buy', 'seven',
 'petrochemical', 'plants',
 'along', 'the', 'texas', 'gulf', 'coast', '.', 'the', 'plants', 'produce',
 'only', 'basic',
 'commodity', 'petrochemicals', 'that', 'are', 'the', 'building', 'blocks',
 'of', 'specialty',
 'products', '.', '"', 'this', 'kind', 'of', 'commodity', 'chemical',
 'business', 'will', 'never',
 'be', 'a', 'glamorous', ',,', 'high', '-', 'margin', 'business', ',,"', 'cain',
 'said', ',,',
 'adding', 'that', 'demand', 'is', 'expected', 'to', 'grow', 'by', 'about',
 'three', 'pct',
 'annually', '.', 'garo', 'armen', ',,', 'an', 'analyst', 'with', 'dean',
 'witter', 'reynolds', ',,',
 'said', 'chemical', 'makers', 'have', 'also', 'benefitted', 'by',
 'increasing', 'demand', 'for',
 'plastics', 'as', 'prices', 'become', 'more', 'competitive', 'with',
 'aluminum', ',,', 'wood',
 'and', 'steel', 'products', '.', 'armen', 'estimated', 'the', 'upturn', 'in',
 'the', 'chemical',
 'business', 'could', 'last', 'as', 'long', 'as', 'four', 'or', 'five',
 'years', ',,', 'provided',
 'the', 'u', '.', 's', '.', 'economy', 'continues', 'its', 'modest', 'rate',
 'of', 'growth', '.',


```

'<END>'],
['<START>', 'turkey', 'calls', 'for', 'dialogue', 'to', 'solve', 'dispute',
'turkey', 'said',
'today', 'its', 'disputes', 'with', 'greece', ',', 'including', 'rights',
'on', 'the',
'continental', 'shelf', 'in', 'the', 'aegean', 'sea', ',', 'should', 'be',
'solved', 'through',
'negotiations', '.', 'a', 'foreign', 'ministry', 'statement', 'said', 'the',
'latest', 'crisis',
'between', 'the', 'two', 'nato', 'members', 'stemmed', 'from', 'the',
'continental', 'shelf',
'dispute', 'and', 'an', 'agreement', 'on', 'this', 'issue', 'would', 'effect',
'the', 'security',
',', 'economy', 'and', 'other', 'rights', 'of', 'both', 'countries', '.', '"',
'as', 'the',
'issue', 'is', 'basically', 'political', ',', 'a', 'solution', 'can', 'only',
'be', 'found', 'by',
'bilateral', 'negotiations', ',',"', 'the', 'statement', 'said', '.', 'greece',
'has', 'repeatedly',
'said', 'the', 'issue', 'was', 'legal', 'and', 'could', 'be', 'solved', 'at',
'the',
'international', 'court', 'of', 'justice', '.', 'the', 'two', 'countries',
'approached', 'armed',
'confrontation', 'last', 'month', 'after', 'greece', 'announced', 'it',
'planned', 'oil',
'exploration', 'work', 'in', 'the', 'aegean', 'and', 'turkey', 'said', 'it',
'would', 'also',
'search', 'for', 'oil', '.', 'a', 'face', '-', 'off', 'was', 'averted',
'when', 'turkey',
'confined', 'its', 'research', 'to', 'territorial', 'waters', '.', '"',
'the', 'latest',
'crises', 'created', 'an', 'historic', 'opportunity', 'to', 'solve', 'the',
'disputes', 'between',
'the', 'two', 'countries', ',',"', 'the', 'foreign', 'ministry', 'statement',
'said', '.', 'turkey',
'"', 's', 'ambassador', 'in', 'athens', ',', 'nazmi', 'akiman', ',', 'was',
'due', 'to', 'meet',
'prime', 'minister', 'andreas', 'papandreou', 'today', 'for', 'the', 'greek',
'reply', 'to', 'a',
'message', 'sent', 'last', 'week', 'by', 'turkish', 'prime', 'minister',
'turgut', 'ozal', '.',
'the', 'contents', 'of', 'the', 'message', 'were', 'not', 'disclosed', '.',
'<END>']]

```

2.2.3 Question 1.1: Implement `distinct_words` [code] (5 points)

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](#).

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

```
[4]: 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): list of distinct words across the
            → corpus, sorted (using python 'sorted' function)
            num_corpus_words (integer): number of distinct words across the
            → corpus
    """
    corpus_words = []
    num_corpus_words = -1

    # -----
    # Write your implementation here.
    corpus_words = sorted(list(set([word for document in corpus for word in
    → document])))
    num_corpus_words = len(corpus_words)
    # -----

    return corpus_words, num_corpus_words

[5]: # -----
    # Run this sanity check
    # Note that this not an exhaustive check for correctness.
    # -----

    # Define toy corpus
    test_corpus = ["START All that glitters isn't gold END".split(" "), "START
    → All's well that ends well END".split(" ")]
    test_corpus_words, num_corpus_words = distinct_words(test_corpus)

    # Correct answers
    ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that",
    → "gold", "All's", "glitters", "isn't", "well", "END"])))
    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!

2.2.4 Question 1.2: Implement compute_co_occurrence_matrix [code] (15 points)

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.

```

[6]: 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 (numpy matrix of shape (number of corpus words, number of corpus
    →words)):
        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 function.
        word2Ind (dict): dictionary that maps word to index (i.e. row/
    →column number) for matrix M.
    """
    words, num_words = distinct_words(corpus)
    M = None
    word2Ind = {}

```

```

# -----
# Write your implementation here.
M = np.zeros((num_words, num_words))
i = 0
for word in words:
    word2Ind[word] = i
    i += 1
for document in corpus:
    for i in range(len(document)):
        current_word = document[i]
        current_word_index = word2Ind[current_word]
        for j in range(1, window_size+1):
            if i-j >= 0:
                co_word_1 = document[i-j]
                co_word_1_index = word2Ind[co_word_1]
                M[current_word_index][co_word_1_index] += 1
            if i+j <= len(document) - 1:
                co_word_2 = document[i+j]
                co_word_2_index = word2Ind[co_word_2]
                M[current_word_index][co_word_2_index] += 1

# -----

return M, word2Ind

```

```

[7]: # -----
# 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 = ["START All that glitters isn't gold END".split(" "), "START_
↳All's well that ends well END".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., 1., 0., 0., 0., 0., 1., 0.,],
     [0., 0., 0., 1., 0., 0., 0., 0., 0., 1.,],
     [0., 0., 0., 0., 0., 0., 1., 0., 0., 1.,],
     [1., 1., 0., 0., 0., 0., 0., 0., 0., 0.,],
     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.,],
     [0., 0., 1., 0., 0., 0., 0., 1., 0., 0.,],
     [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,],
     [1., 0., 0., 0., 1., 1., 0., 0., 0., 1.,],
     [0., 1., 1., 0., 1., 0., 0., 0., 1., 0.,]]

```

```

)
word2Ind_ans = {'All': 0, "All's": 1, 'END': 2, 'START': 3, 'ends': 4,
→'glitters': 5, 'gold': 6, "isn't": 7, 'that': 8, 'well': 9}

# 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, w1, w2,
→student, correct))

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

```

Passed All Tests!

2.2.5 Question 1.3: Implement `reduce_to_k_dim` [code] (10 point)

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](#).

```
[8]: 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 corpus words, number of corpus
    →words)): co-occurrence matrix of word counts
        k (int): embedding size of each word after dimension reduction
    Return:
        M_reduced (numpy matrix of shape (number of corpus words, k)):
    →matrix of k-dimensional 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`
    M_reduced = None
    print("Running Truncated SVD over %i words..." % (M.shape[0]))

    # -----
    # Write your implementation here.
    M_reduced = TruncatedSVD(n_components = k, n_iter = n_iters).
    →fit_transform(M)
    # -----

    print("Done.")
    return M_reduced
```

```
[9]: # -----
# Run this sanity check
# Note that this 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 = ["START All that glitters isn't gold END".split(" "), "START
    →All's well that ends well END".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)

```

Running Truncated SVD over 10 words...
Done.

Passed All Tests!

2.2.6 Question 1.4: Implement plot_embeddings [code] (10 point)

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.

```

[10]: 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 , k)): matrix of k-dimensional word embeddings
        word2Ind (dict): dictionary that maps word to indices for matrix M
        words (list of strings): words whose embeddings we want to_
    →visualize
    """

    # -----
    # Write your implementation here.
    for word in words:
        index = word2Ind[word]
        coordinates = M_reduced[index]
        x = coordinates[0]
        y = coordinates[1]
        plt.scatter(x, y, marker='x', color='red')
        plt.text(x, y, word, fontsize=9)
    # -----

```

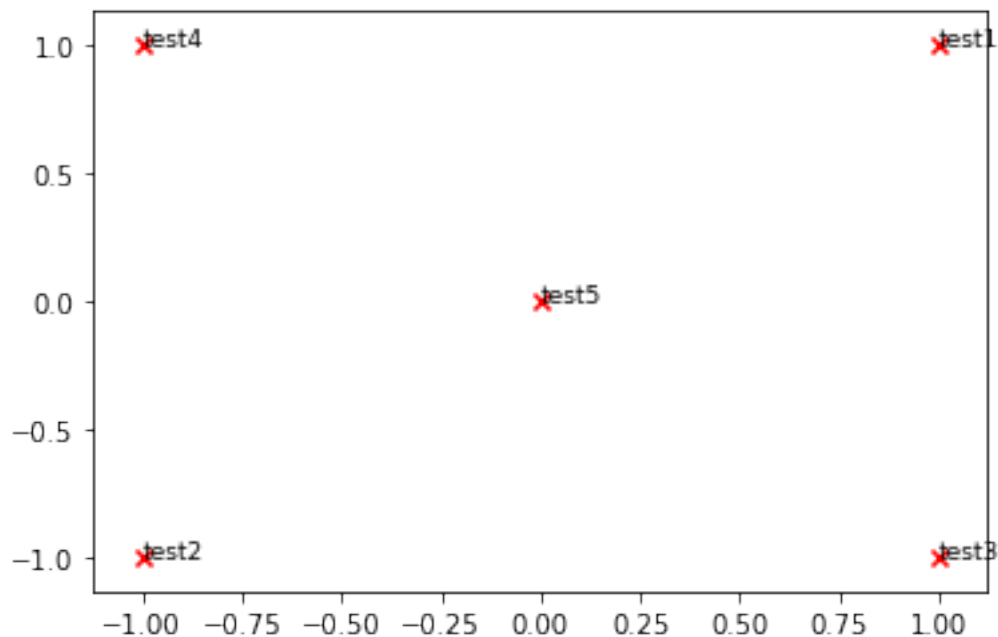
```
[11]: # -----
# Run this sanity check
# Note that this not an exhaustive check for correctness.
# The plot produced should look like the "test solution plot" depicted below.
# -----

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:



Test Plot Solution

2.2.7 Question 1.5: Co-Occurrence Plot Analysis [written] (10 points)

Now we will put together all the parts you have written! We will compute the co-occurrence matrix with fixed window of 4, over the Reuters “crude” corpus. Then we will use TruncatedSVD

to compute 2-dimensional embeddings of each word. TruncatedSVD returns $U \cdot S$, so we 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. What clusters together in 2-dimensional embedding space? What doesn't cluster together that you might think should have? **Note:** "bpd" stands for "barrels per day" and is a commonly used abbreviation in crude oil topic articles.

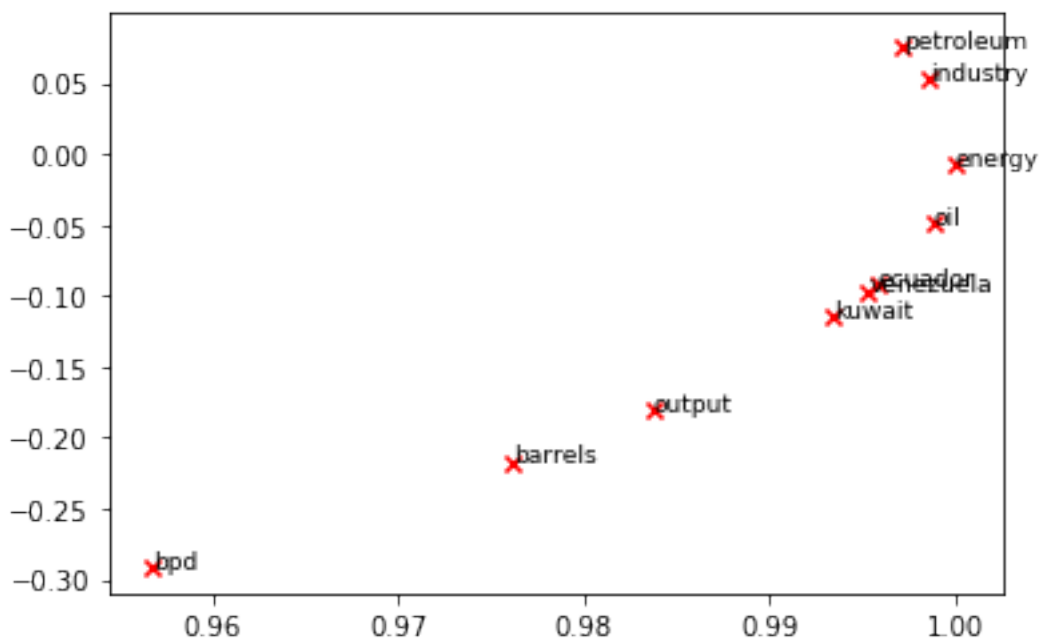
```
[12]: # -----
# Run This Cell to Produce Your Plot
# -----
reuters_corpus = read_corpus()
M_co_occurrence, word2Ind_co_occurrence = \
    compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

# 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 = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil',
        'output', 'petroleum', 'venezuela']
plot_embeddings(M_normalized, word2Ind_co_occurrence, words)
```

Running Truncated SVD over 8185 words...

Done.



Write your answer here. From my perspective, there are mainly five clusters: * ecuador, kuwait, venezuela: Three countries are clustered together, they are all oil-producing countries.

- energy, oil: Oil can generate energy, it is not surprising that they are clustered together.
- petroleum, industry: The petroleum industry is a terminology, so it is reasonable that “petroleum” and “industry” are clustered together.
- barrels, output: The output of the petroleum industry is measured by barrels, which explains why they are clustered together.
- bpd: Barrels per day is left aside by our model.

I think “bpd”, “barrels”, “output” should be clustered together, since barrels per day measures the output of the petroleum industry, however, “bpd” is not clustered with “barrels” and “output”. I think this happens because we either use “bpd” or “barrels per day” to describe the amount of output but not both of them, so “bpd” and “barrels” might have a relationship that is somehow mutually-exclusive.

2.3 Part 2: Prediction-Based Word Vectors (50 points)

More recently prediction-based word vectors have come into fashion, e.g. word2vec. Here, we shall explore the embeddings produced by word2vec. If you’re feeling adventurous, challenge yourself and try reading the [original paper](#).

Then run the following cells to load the word2vec vectors into memory. **Note:** This might take several minutes.

```
[13]: def load_word2vec():  
    """ Load Word2Vec Vectors  
    Return:  
        wv_from_bin: All 3 million embeddings, each length 300  
    """  
    import gensim.downloader as api  
    wv_from_bin = api.load("word2vec-google-news-300")  
    vocab = list(wv_from_bin.vocab.keys())  
    print("Loaded vocab size %i" % len(vocab))  
    return wv_from_bin
```

```
[14]: # -----  
# Run Cell to Load Word Vectors  
# Note: This may take several minutes  
# -----  
wv_from_bin = load_word2vec()
```

Loaded vocab size 3000000

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.

2.3.1 Reducing dimensionality of Word2Vec Word Embeddings

Let's directly compare the word2vec embeddings to those of the co-occurrence matrix. Run the following cells to:

1. Put the 3 million word2vec vectors into a matrix M
2. Run `reduce_to_k_dim` (your Truncated SVD function) to reduce the vectors from 300-dimensional to 2-dimensional.

```
[15]: def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd',  
    → 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum',  
    → 'venezuela']):  
    """ Put the word2vec vectors into a matrix M.  
    Param:  
        wv_from_bin: KeyedVectors object; the 3 million word2vec vectors  
    → loaded from file  
    Return:  
        M: numpy matrix shape (num words, 300) containing the vectors  
        word2Ind: dictionary mapping each word to its row number in M  
    """  
    import random  
    words = list(wv_from_bin.vocab.keys())  
    print("Shuffling words ...")  
    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.word_vec(w))  
            word2Ind[w] = curInd  
            curInd += 1  
        except KeyError:  
            continue  
    for w in required_words:  
        try:  
            M.append(wv_from_bin.word_vec(w))  
            word2Ind[w] = curInd  
            curInd += 1  
        except KeyError:  
            continue
```

```

M = np.stack(M)
print("Done.")
return M, word2Ind

```

```

[16]: # -----
# Run Cell to Reduce 300-Dimensional Word Embeddings to k Dimensions
# Note: This may take several minutes
# -----
M, word2Ind = get_matrix_of_vectors(wv_from_bin)
M_reduced = reduce_to_k_dim(M, k=2)

```

Shuffling words ...

Putting 10000 words into word2Ind and matrix M...

Done.

Running Truncated SVD over 10010 words...

Done.

2.3.2 Question 2.1: Word2Vec Plot Analysis [written] (5 points)

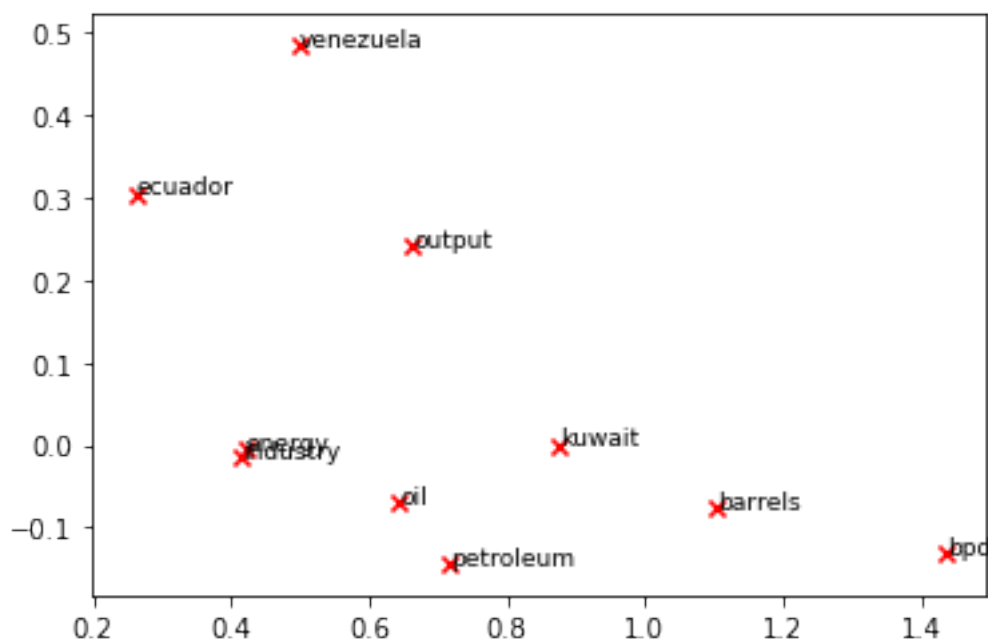
Run the cell below to plot the 2D word2vec embeddings for ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum', 'venezuela'].

What clusters together in 2-dimensional embedding space? What doesn't cluster together that you might think should have? How is the plot different from the one generated earlier from the co-occurrence matrix?

```

[17]: words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output',
    → 'petroleum', 'venezuela']
plot_embeddings(M_reduced, word2Ind, words)

```



Write your answer here. What clusters together in 2-dimensional embedding space?

There is only one obvious cluster, which is, * energy, industry

Other word vectors are scattered around.

What doesn't cluster together that you might think should have?

Three countries, Ecuador, Venezuela, Kuwait are no longer clustered together which I think should be clustered together. Also, "barrel" and "output" are not clustered together, but they always co-occur to describe the measure of the output of the oil industry.

How is the plot different from the one generated earlier from the co-occurrence matrix?

In this plot, word vectors are scattered around and there is only one obvious cluster; but in the plot generated from the co-occurrence matrix, there are multiple clusters.

2.3.3 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:

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

2.3.4 Question 2.2: Polysemous Words (5 points) [code + written]

Find a [polysemous](#) word (for example, "leaves" or "scoop") such that the top-10 most similar words (according to cosine similarity) contains related words from *both* meanings. For example, "leaves" has both "vanishes" and "stalks" in the top 10, and "scoop" has both "handed_waffle_cone" and "lowdown". You will probably need to try several polysemous words before you find one. Please state the polysemous word you discover and the multiple meanings that occur in the top 10. Why do you think many of the polysemous words you tried didn't work?

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](#).

```
[18]: # -----  
# Write your polysemous word exploration code here.  
wv_from_bin.most_similar("string")  
# -----
```

```
[18]: [('spate', 0.6308412551879883),  
      ('slew', 0.5636824369430542),  
      ('rash', 0.5049306750297546),
```

```
('litany', 0.4844335615634918),
('String', 0.4780084788799286),
('flurry', 0.4644429385662079),
('strung_together', 0.4589213728904724),
('stringing_together', 0.4524788558483124),
('bevy', 0.4505397081375122),
('trio', 0.4444049596786499)]
```

Write your answer here.

- The polysemous word I discovered: string

The top 10 similar words of “string” are: spate, slew, rash, litany, String, flurry, strung_together, stringing_together, bevy, trio

The meanings they reflect: 1. a cord usually used to bind, fasten, or tie; 2. hang (something) so that it stretches in a long line; 3. a data type in Computer Science

- Some polysemous words I tried that didn’t work: bank, spring, state

The reason I think some polysemous words do not work is that they are always used to express a specific meaning. For example, “bank” always used to refer to the financial institution instead of the slopping land near river. Or in other words, “bank” are more frequently used with the meaning, a financial institution.

2.3.5 Question 2.3: Synonyms & Antonyms (5 points) [code + written]

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 $\text{Cosine Distance}(w_1, w_3) < \text{Cosine Distance}(w_1, w_2)$. For example, $w_1 = \text{“happy”}$ is closer to $w_3 = \text{“sad”}$ than to $w_2 = \text{“cheerful”}$.

Once you have found your example, please give a possible explanation for why this counter-intuitive result may have happened.

You should use 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.

```
[19]: # -----
# Write your synonym & antonym exploration code here.

w1 = "good"
w2 = "excellent"
w3 = "bad"
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))

# -----
```

Synonyms good, excellent have cosine distance: 0.35570716857910156
Antonyms good, bad have cosine distance: 0.28099489212036133

Write your answer here. Although “good” and “excellent” have similar meaning, they are both adjectives to describe something that is positive, they are sometimes not interchangeable, for example, when we say “do not wait for good things happen to you, you need to walk toward happiness”, it is strange to say “do not wait for excellent things happen to you, you need to walk toward happiness”. However, since “good” and “bad” are antonyms, they are often used with the same grammatical structure and similar context, the only difference is their meanings are opposite, in other words, they are interchangeable if we want to express the opposite meaning.

2.3.6 Solving Analogies with Word Vectors

Word2Vec vectors have been shown to *sometimes* exhibit the ability to solve analogies.

As an example, for the analogy “man : king :: woman : x”, what is x?

In the cell below, we show you how to use word vectors to find x. The `most_similar` function finds words that are most similar to the words in the positive list and most dissimilar from the words in the negative list. The answer to the analogy will be the word ranked most similar (largest numerical value).

Note: Further Documentation on the `most_similar` function can be found within the [GenSim documentation](#).

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

```
[('queen', 0.7118192911148071),
 ('monarch', 0.6189674139022827),
 ('princess', 0.5902431607246399),
 ('crown_prince', 0.5499460697174072),
 ('prince', 0.5377321243286133),
 ('kings', 0.5236844420433044),
 ('Queen_Consort', 0.5235945582389832),
 ('queens', 0.5181134343147278),
 ('sultan', 0.5098593235015869),
 ('monarchy', 0.5087411999702454)]
```

2.3.7 Question 2.4: Finding Analogies [code + written] (10 Points)

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!

```
[21]: # -----
# Write your analogy exploration code here.

pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'father'],
    ↪negative=['man']))
```

```
# -----
```

```
[('mother', 0.8462507128715515),  
 ('daughter', 0.7899606227874756),  
 ('husband', 0.7560455799102783),  
 ('son', 0.7279756665229797),  
 ('eldest_daughter', 0.7120418548583984),  
 ('niece', 0.7096832990646362),  
 ('aunt', 0.6960804462432861),  
 ('grandmother', 0.6897341012954712),  
 ('sister', 0.6895190477371216),  
 ('daughters', 0.6731119751930237)]
```

Write your answer here. The analogy is “man : father :: woman : mother”.

2.3.8 Question 2.5: Incorrect Analogy [code + written] (5 point)

Find an 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.

```
[22]: # -----  
# Write your incorrect analogy exploration code here.  
  
pprint.pprint(wv_from_bin.most_similar(positive=['female', 'grandfather'],  
→negative=['male']))  
  
# -----
```

```
[('grandson', 0.7096990942955017),  
 ('granddaughter', 0.670899510383606),  
 ('father', 0.6570508480072021),  
 ('grandmother', 0.6495229005813599),  
 ('uncle', 0.6471563577651978),  
 ('greatgrandfather', 0.6320098638534546),  
 ('paternal_grandfather', 0.6265338659286499),  
 ('maternal_grandfather', 0.61909019947052),  
 ('niece', 0.5914374589920044),  
 ('aunt', 0.589950680732727)]
```

Write your answer here. The analogy should be “male : grandfather :: female : grandmother”, but the code gives “male : grandfather :: female : grandson”, “granson” is incorrect.

2.3.9 Question 2.6: Guided Analysis of Bias in Word Vectors [written] (5 point)

It’s important to be cognizant of the biases (gender, race, sexual orientation etc.) implicit to our word embeddings.

Run the cell below, to examine (a) which terms are most similar to “woman” and “worker” and most dissimilar to “man”, and (b) which terms are most similar to “man” and “worker” and most dissimilar to “woman”. What do you find in the top 10?

```
[23]: # Run this cell
# Here `positive` indicates the list of words to be similar to and `negative`
      ↳ indicates the list of words to be
# most dissimilar from.
pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'worker'],
      ↳ negative=['man']))
print()
pprint.pprint(wv_from_bin.most_similar(positive=['man', 'worker'],
      ↳ negative=['woman']))
```

```
[('workers', 0.6582455635070801),
 ('employee', 0.5805294513702393),
 ('nurse', 0.5249922275543213),
 ('receptionist', 0.5142490267753601),
 ('migrant_worker', 0.5001609921455383),
 ('Worker', 0.4979270100593567),
 ('housewife', 0.48609834909439087),
 ('registered_nurse', 0.4846190810203552),
 ('laborer', 0.48437264561653137),
 ('coworker', 0.48212409019470215)]
```

```
[('workers', 0.5590359568595886),
 ('laborer', 0.54481041431427),
 ('foreman', 0.5192232131958008),
 ('Worker', 0.5161596536636353),
 ('employee', 0.5094279646873474),
 ('electrician', 0.49481213092803955),
 ('janitor', 0.48718899488449097),
 ('bricklayer', 0.4825313687324524),
 ('carpenter', 0.47499001026153564),
 ('workman', 0.4642517566680908)]
```

Write your answer here. which terms are most similar to “woman” and “worker” and most dissimilar to “man”

- “nurse”, “receptionist”, “housewife”, are some terms are most similar to “woman” and “worker” and most dissimilar to “man”.

which terms are most similar to “man” and “worker” and most dissimilar to “woman”

- “laborer”, “foreman”, “electrician”, “bricklayer”, “carpenter”, are some terms are most similar to “man” and “worker” and most dissimilar to “woman”.

In general, I find that some terms of occupation are highly gender-biased, for example, nurse, receptionist, housewife are biased towards female since the majority of workers in these areas

are women; laborer, foreman, electrician, bricklayer, carpenter are biased towards male since the majority of workers in these areas are men.

2.3.10 Question 2.7: Independent Analysis of Bias in Word Vectors [code + written] (10 points)

Use the `most_similar` function to find another case where some bias is exhibited by the vectors. Please briefly explain the example of bias that you discover.

```
[24]: # -----  
# Write your bias exploration code here.  
  
pprint.pprint(wv_from_bin.most_similar(positive=['Chinese', 'worker'],  
→negative=['American']))  
print()  
pprint.pprint(wv_from_bin.most_similar(positive=['American', 'worker'],  
→negative=['Chinese']))  
  
# -----
```

```
[('migrant_worker', 0.5319787263870239),  
 ('workers', 0.5307591557502747),  
 ('ayi', 0.5247678756713867),  
 ('surnamed_Song', 0.5071259140968323),  
 ('surnamed_Xiao', 0.5010257959365845),  
 ('surnamed_Yang', 0.5005248785018921),  
 ('surnamed_Huang', 0.500294029712677),  
 ('surnamed_Ma', 0.4971430003643036),  
 ('surnamed_Li', 0.49599695205688477),  
 ('surnamed_Zhu', 0.49494844675064087)]  
  
[('employee', 0.4758417010307312),  
 ('workers', 0.45908257365226746),  
 ('forklift_operator', 0.39466822147369385),  
 ('laborer', 0.3933212161064148),  
 ('janitor', 0.3807808756828308),  
 ('electrician', 0.37922415137290955),  
 ('pipe_fitter', 0.37634676694869995),  
 ('electrician_apprentice', 0.3761669993400574),  
 ('X_ray_technician', 0.37583503127098083),  
 ('technician', 0.37248194217681885)]
```

Write your answer here. Chinese workers' surnames are often biased towards surnames such as Song, Xiao, Yang, Huang, Ma, Li, and Zhu. Also, they are considered as migrant workers in the U.S.

American workers are often refer to employees or laborers, furthermore, electrician, forklift operator, pipe fitter, technician, are some of representative occupations of American workers.

2.3.11 Question 2.8: Thinking About Bias [written] (5 point)

What might be the cause of these biases in the word vectors?

Write your answer here. As Prof. Bhat mentioned in class, “words that occur in similar contexts tend to have similar meanings”. Because some words have similar meanings, they often occur together, therefore we see a “bias” that these words are often correlated.

Also, since our society have some biases, such as bias towards gender, that the majority of some occupations are from one gender instead of both, these biases in our society are reflected in our language, therefore, there are biases in the word vectors.

3 Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.
2. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
3. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.
4. Once you’ve rerun everything, select File -> Download as -> PDF via LaTeX, and make sure to place it in the same directory as the assignment.
5. Look at the PDF file and make sure all your solutions are there, displayed correctly.
6. Then, compress (tar) the complete directory of the assignment.
7. Submit your tar file on Compass.