

▼ XCS224N Assignment 1: Exploring Word Vectors (24 Points)

Before you start, make sure you read "XCS224N HW1 - Handout".

```
1  # All Import Statements Defined Here
2  # Note: Do not add to this list.
3  # All the dependencies you need can be installed by running this cell.
4  # Throughout this notebook you can run a cell by hitting CTRL+RETURN or the Play button/icon at left
5  # -----
6
7  import sys
8  assert sys.version_info[0]==3
9  assert sys.version_info[1] >= 5
10
11  from gensim.models import KeyedVectors
12  from gensim.test.utils import datapath
13  import pprint
14  import matplotlib.pyplot as plt
15  plt.rcParams['figure.figsize'] = [10, 5]
16  import nltk
17  nltk.download('reuters')
18  from nltk.corpus import reuters
19  import numpy as np
20  import random
21  import scipy as sp
22  from sklearn.decomposition import TruncatedSVD
23  from sklearn.decomposition import PCA
24
25  START_TOKEN = '<START>'
26  END_TOKEN = '<END>'
27
28  np.random.seed(0)
29  random.seed(0)
30  # -----
```

```
📄 [nltk_data] Downloading package reuters to /root/nltk_data...
[nltk_data]   Package reuters is already up-to-date!
```

▼ Your Name: Wei Wu

Your SCPD XID Number: X484700

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

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*.

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 (10 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](#) or [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.

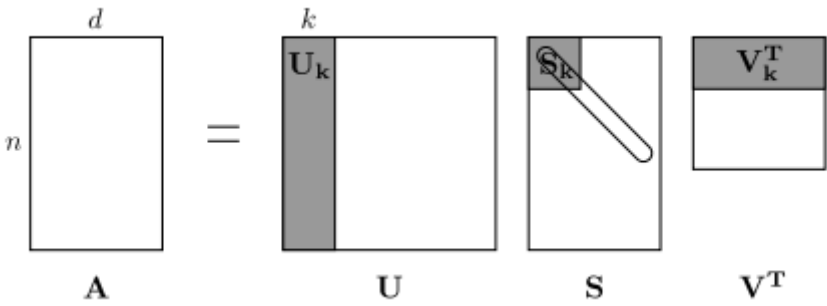
Example: Co-Occurrence with Fixed Window of n=1:

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	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 these 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 CTRL-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 perform any other kind of pre-processing.

```
1 def read_corpus(category="crude"):  
2     """ Read files from the specified Reuter's category.  
3     Params:  
4         category (string): category name  
5     Return:  
6         list of lists, with words from each of the processed files  
7     """  
8     files = reuters.fileids(category)  
9     return [[START_TOKEN] + [w.lower() for w in list(reuters.words(f))] + [END_TOKEN] for f in files]  
10
```

Let's have a look what these documents are like.... (run the cell below)

```
1 reuters_corpus = read_corpus()  
2 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',
```

▼ Question 1.1: Implement `distinct_words` [code] (2 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.

```
1 def distinct_words(corpus):
2     """ Determine a list of distinct words for the corpus.
3         Params:
4             corpus (list of list of strings): corpus of documents
5         Return:
6             corpus_words (list of strings): list of distinct words across the corpus, sorted (using python 'sorted'
7             num_corpus_words (integer): number of distinct words across the corpus
8     """
9     corpus_words = []
10    num_corpus_words = -1
11
12    ### SOLUTION BEGIN
13    corpus_words = sorted(list({word for doc in corpus for word in doc}))
14    num_corpus_words = len(corpus_words)
15
16    ### SOLUTION END
17
18    return corpus_words, num_corpus_words

```

```
1 # -----
2 # Run this sanity check
3 # Note that this not an exhaustive check for correctness.
4 "
```

```

4  # -----
5
6  # Define toy corpus
7  test_corpus = ["START All that glitters isn't gold END".split(" "), "START All's well that ends well END".split(" ")
8  test_corpus_words, num_corpus_words = distinct_words(test_corpus)
9
10 # Correct answers
11 ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well"
12 ans_num_corpus_words = len(ans_test_corpus_words)
13
14 # Test correct number of words
15 assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of distinct words. Correct: {}. Yours: {}".format
16
17 # Test correct words
18 assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.\nCorrect: {}\nYours:   {}".format(str(a
19
20 # Print Success
21 print ("-" * 80)
22 print("Passed All Tests!")
23 print ("-" * 80)

```

```

➦ -----
Passed All Tests!
-----

```

▼ Question 1.2: Implement `compute_co_occurrence_matrix` [code] (3 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. If you're not familiar with NumPy, there's a NumPy tutorial in the second half of this cs231n [Python NumPy tutorial](#).

```

1  from itertools import chain # noticed that I import itertools here
2
3  def compute_co_occurrence_matrix(corpus, window_size=4):
4      """ Compute co-occurrence matrix for the given corpus and window_size (default of 4).
5
6          Note: Each word in a document should be at the center of a window. Words near edges will have a smaller
7              number of co-occurring words.
8
9          For example, if we take the document "START All that glitters is not gold END" with window size of 4,
10             "All" will co-occur with "START", "that", "glitters", "is", and "not".
11
12         Params:
13             corpus (list of list of strings): corpus of documents
14             window_size (int): size of context window
15         Return:
16             M (numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus):
17                 Co-occurrence matrix of word counts.
18                 The ordering of the words in the rows/columns should be the same as the ordering of the words given l
19             word2Ind (dict): dictionary that maps word to index (i.e. row/column number) for matrix M.
20         """
21         words, num_words = distinct_words(corpus)
22         M = None
23         word2Ind = {}
24
25     ### SOLUTION BEGIN
26     M = np.zeros((num_words,num_words))
27     for idx, word in enumerate(words):
28         word2Ind[word] = idx
29
30     for doc in corpus:
31         for i, word in enumerate(doc):
32             c_idx = word2Ind[word]
33             # i-n to i-1 and i+1 to i+n, truncated
34             for j in chain(range(max(i - window_size,0), i), range(i + 1, min(i + window_size + 1, len(doc)))):
35                 o_idx = word2Ind[doc[j]]
36                 M[c_idx, o_idx] += 1
37     ### SOLUTION END
38
39     return M, word2Ind

```

```

1  # -----
2  # Run this sanity check
3  # Note that this is not an exhaustive check for correctness.
4  # -----
5
6  # Define toy corpus and get student's co-occurrence matrix
7  test_corpus = ["START All that glitters isn't gold END".split(" "), "START All's well that ends well END".split(" ")
8  test_corpus_words, num_corpus_words = distinct_words(test_corpus)
9
10 # Correct answers
11 ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well"
12 ans_num_corpus_words = len(ans_test_corpus_words)
13
14 # Test correct number of words
15 assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of distinct words. Correct: {}. Yours: {}".format
16
17 # Test correct words
18 assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.\nCorrect: {}\nYours:   {}".format(str(a
19
20 # Print Success
21 print ("-" * 80)
22 print("Passed All Tests!")
23 print ("-" * 80)

```

```

8 M_test, word2Ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)
9
10 # Correct M and word2Ind
11 M_test_ans = np.array(
12     [[0., 0., 0., 1., 0., 0., 0., 0., 1., 0.,],
13     [0., 0., 0., 1., 0., 0., 0., 0., 0., 1.,],
14     [0., 0., 0., 0., 0., 0., 1., 0., 0., 1.,],
15     [1., 1., 0., 0., 0., 0., 0., 0., 0., 0.,],
16     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.,],
17     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0.,],
18     [0., 0., 1., 0., 0., 0., 0., 1., 0., 0.,],
19     [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,],
20     [1., 0., 0., 0., 1., 1., 0., 0., 0., 1.,],
21     [0., 1., 1., 0., 1., 0., 0., 0., 1., 0.,]]
22 )
23 word2Ind_ans = {'All': 0, "All's": 1, 'END': 2, 'START': 3, 'ends': 4, 'glitters': 5, 'gold': 6, "isn't": 7, 'that':
24
25 # Test correct word2Ind
26 assert (word2Ind_ans == word2Ind_test), "Your word2Ind is incorrect:\nCorrect: {}\nYours: {}".format(word2Ind_ans, word2Ind_test)
27
28 # Test correct M shape
29 assert (M_test.shape == M_test_ans.shape), "M matrix has incorrect shape.\nCorrect: {}\nYours: {}".format(M_test.shape, M_test_ans.shape)
30
31 # Test correct M values
32 for w1 in word2Ind_ans.keys():
33     idx1 = word2Ind_ans[w1]
34     for w2 in word2Ind_ans.keys():
35         idx2 = word2Ind_ans[w2]
36         student = M_test[idx1, idx2]
37         correct = M_test_ans[idx1, idx2]
38         if student != correct:
39             print("Correct M:")
40             print(M_test_ans)
41             print("Your M: ")
42             print(M_test)
43             raise AssertionError("Incorrect count at index ({}, {})=({}, {}) in matrix M. Yours has {} but should have {}".format(idx1, idx2, correct, student, student, correct))
44
45 # Print Success
46 print("-" * 80)
47 print("Passed All Tests!")
48 print("-" * 80)

```



Passed All Tests!

▼ Question 1.3: Implement `reduce_to_k_dim` [code] (1 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](http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html).

```

1 def reduce_to_k_dim(M, k=2):
2     """ Reduce a co-occurrence count matrix of dimensionality (num_corpus_words, num_corpus_words)
3         to a matrix of dimensionality (num_corpus_words, k) using the following SVD function from Scikit-Learn:
4         - http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html
5
6     Params:
7         M (numpy matrix of shape (number of unique words in the corpus , number of unique words in the corpus)):
8         k (int): embedding size of each word after dimension reduction
9     Return:
10         M_reduced (numpy matrix of shape (number of corpus words, k)): matrix of k-dimensional word embeddings.
11         In terms of the SVD from math class, this actually returns U * S
12
13     """
14     n_iters = 10 # Use this parameter in your call to `TruncatedSVD`
15     M_reduced = None
16     print("Running Truncated SVD over %i words..." % (M.shape[0]))
17
18     ### SOLUTION BEGIN
19     svd = TruncatedSVD(n_components = k, n_iter = n_iters, random_state= 42)
20     M_reduced = svd.fit_transform(M)
21     ### SOLUTION END
22
23     print("Done.")
24     return M_reduced

```

1 # -----


```

2 # Run this sanity check
3 # Note that this not an exhaustive check for correctness
4 # In fact we only check that your M_reduced has the right dimensions.
5 # -----
6
7 # Define toy corpus and run student code
8 test_corpus = ["START All that glitters isn't gold END".split(" "), "START All's well that ends well END".split(" ")
9 M_test, word2Ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)
10 M_test_reduced = reduce_to_k_dim(M_test, k=2)
11
12 # Test proper dimensions
13 assert (M_test_reduced.shape[0] == 10), "M_reduced has {} rows; should have {}".format(M_test_reduced.shape[0], 10)
14 assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have {}".format(M_test_reduced.shape[1], 2)
15
16 # Print Success
17 print ("-" * 80)
18 print("Passed All Tests!")
19 print ("-" * 80)

```

➤ Running Truncated SVD over 10 words...
Done.

Passed All Tests!

▼ Question 1.4: Implement `plot_embeddings` [code] (1 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.

```

1 def plot_embeddings(M_reduced, word2Ind, words):
2     """ Plot in a scatterplot the embeddings of the words specified in the list "words".
3         NOTE: do not plot all the words listed in M_reduced / word2Ind.
4         Include a label next to each point.
5
6         Params:
7             M_reduced (numpy matrix of shape (number of unique words in the corpus , k)): matrix of k-dimensioal word
8             word2Ind (dict): dictionary that maps word to indices for matrix M
9             words (list of strings): words whose embeddings we want to visualize
10    """
11
12    ### SOLUTION BEGIN
13    words_idx = [word2Ind[word] for word in words]
14    x_coords = [M_reduced[idx,0] for idx in words_idx]
15    y_coords = [M_reduced[idx,1] for idx in words_idx]
16
17    for i,word in enumerate(words):
18        x = x_coords[i]
19        y = y_coords[i]
20        plt.scatter(x, y, marker='x', color='red')
21        plt.annotate(word, (x, y))
22
23    plt.show()
24
25    ### SOLUTION END

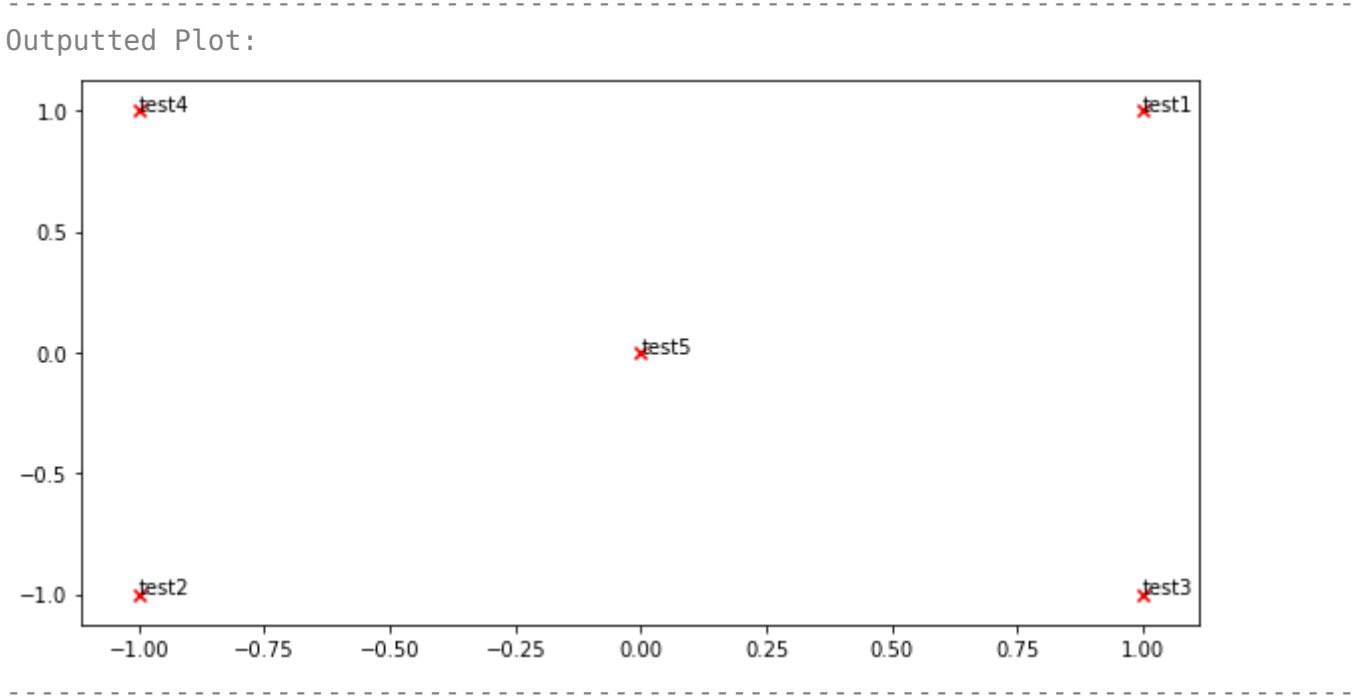
```

```

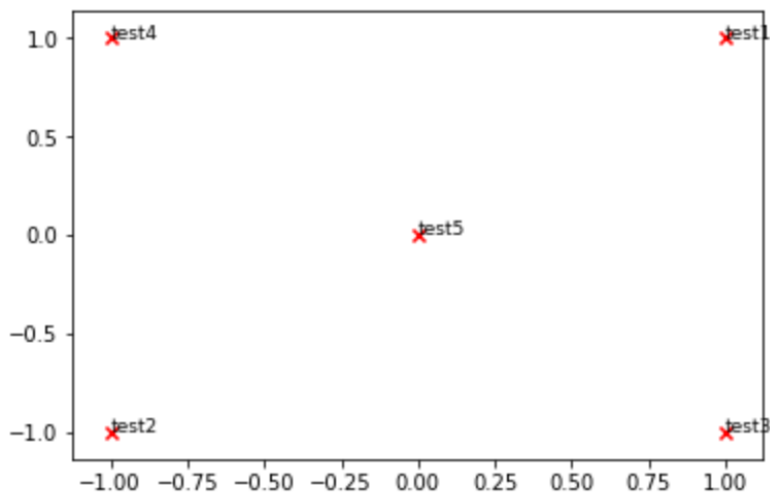
1 # -----
2 # Run this sanity check
3 # Note that this not an exhaustive check for correctness.
4 # The plot produced should look like the "test solution plot" depicted below.
5 # -----
6
7 print ("-" * 80)
8 print ("Outputted Plot:")
9
10 M_reduced_plot_test = np.array([[1, 1], [-1, -1], [1, -1], [-1, 1], [0, 0]])
11 word2Ind_plot_test = {'test1': 0, 'test2': 1, 'test3': 2, 'test4': 3, 'test5': 4}
12 words = ['test1', 'test2', 'test3', 'test4', 'test5']
13 plot_embeddings(M_reduced_plot_test, word2Ind_plot_test, words)
14
15 print ("-" * 80)

```





Test Plot Solution



Question 1.5: Co-Occurrence Plot Analysis [written] (3 points)

Now we will put together all the parts you have written! We will compute the co-occurrence matrix with fixed window of 5, over the Reuters "crude" corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns U*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](#).

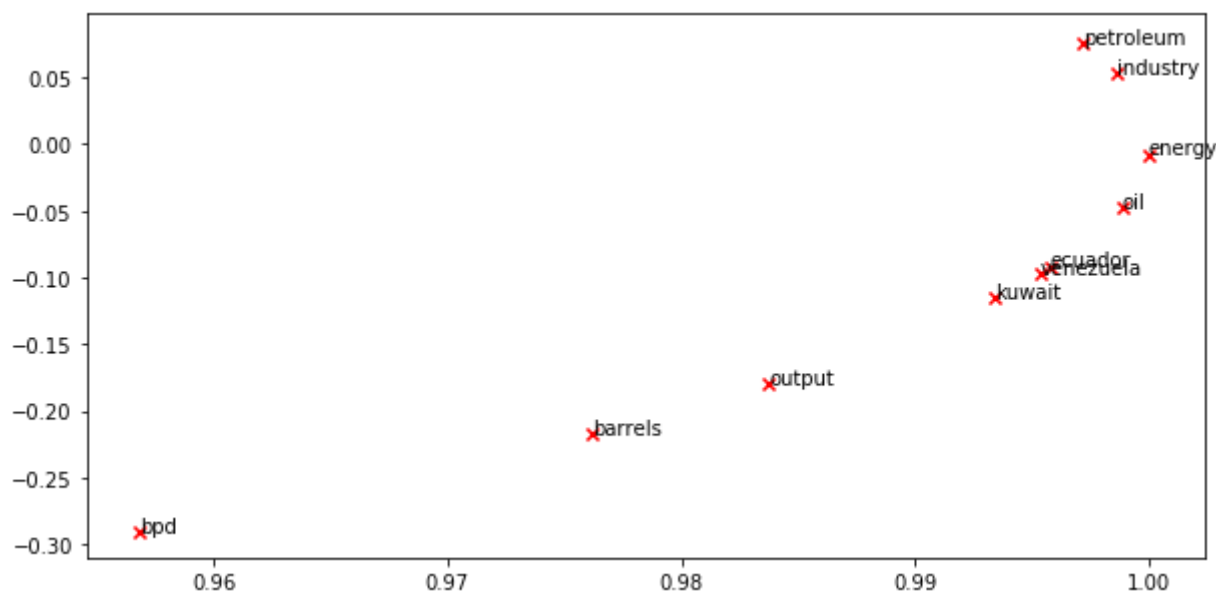
Run the below cell to produce the plot. It'll probably take a few seconds to run.

Written Question: *What clusters together in 2-dimensional embedding space (in the given plot)? 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.*

```
1  # -----
2  # Run This Cell to Produce Your Plot
3  # -----
4  reuters_corpus = read_corpus()
5  M_co_occurrence, word2Ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
6  M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)
7
8  # Rescale (normalize) the rows to make them each of unit-length
9  M_lengths = np.linalg.norm(M_reduced_co_occurrence, axis=1)
10 M_normalized = M_reduced_co_occurrence / M_lengths[:, np.newaxis] # broadcasting
11
12 words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum', 'venezuela']
13 plot_embeddings(M_normalized, word2Ind_co_occurrence, words)
```



Running Truncated SVD over 8185 words...
Done.



Write your answer here.

The words bpd and barrels should be clustered together since they both involve the word barrle. I am not expecting kuwait and ecuador/venezuela to be clustered together since they are not geographically close.

▼ Part 2: Prediction-Based Word Vectors (14 points)

As discussed in class, more recently prediction-based word vectors have come into fashion, e.g. word2vec. Here, we shall explore the embeddings produced by word2vec. Please revisit the class notes and lecture slides for more details on the word2vec algorithm. 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 could take several minutes.**

```
1 def load_word2vec():
2     """ Load Word2Vec Vectors
3     Return:
4         wv_from_bin: 2.5 million of 3 million embeddings, each length 300
5     """
6     import gensim.downloader as api
7     from gensim.models import KeyedVectors
8     # let's load 2.5 million of the 3 million word embeddings so we don't run out of memory on Colab
9     wv_from_bin = KeyedVectors.load_word2vec_format(api.load("word2vec-google-news-300", return_path=True), limit=2500000)
10    vocab = list(wv_from_bin.vocab.keys())
11    print("Loaded vocab size %i" % len(vocab))
12    return wv_from_bin
```

```
1 # -----
2 # Run Cell to Load Word Vectors
3 # Note: This may take several minutes
4 # -----
5 wv_from_bin = load_word2vec()
```

```
➤ /usr/local/lib/python3.6/dist-packages/smart_open/smart_open_lib.py:398: UserWarning: This function is deprecated, use
  'See the migration notes for details: %s' % _MIGRATION_NOTES_URL
Loaded vocab size 2500000
```

▼ 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 2.5 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.

```
1 def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait',
2     """ Put the word2vec vectors into a matrix M.
3     Param:
4         wv_from_bin: KeyedVectors object; the 2.5 million word2vec vectors loaded from file
5     Return:
6         M: numpy matrix shape (num words, 300) containing the vectors
7         word2Ind: dictionary mapping each word to its row number in M
8     """
9     import random
10    words = list(wv_from_bin.vocab.keys())
11    print("Shuffling words ...")
12    random.shuffle(words)
```

```
13 words = words[:10000]
14 print("Putting %i words into word2Ind and matrix M..." % len(words))
15 word2Ind = {}
16 M = []
17 curInd = 0
18 for w in words:
19     try:
20         M.append(wv_from_bin.word_vec(w))
21         word2Ind[w] = curInd
22         curInd += 1
23     except KeyError:
24         continue
25 for w in required_words:
26     try:
27         M.append(wv_from_bin.word_vec(w))
28         word2Ind[w] = curInd
29         curInd += 1
30     except KeyError:
31         continue
32 M = np.stack(M)
33 print("Done.")
34 return M, word2Ind
```

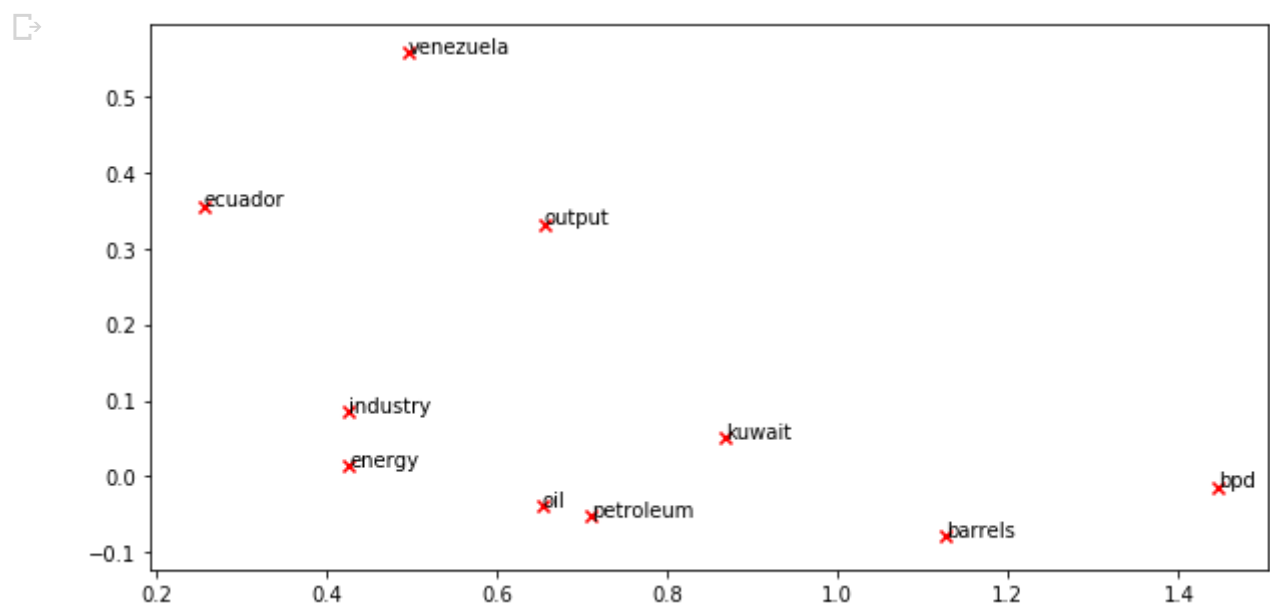
```
1 # -----
2 # Run Cell to Reduce 300-Dimensinal Word Embeddings to k Dimensions
3 # Note: This may take several minutes
4 # -----
5 M, word2Ind = get_matrix_of_vectors(wv_from_bin)
6 M_reduced = reduce_to_k_dim(M, k=2)
```

Shuffling words ...
Putting 10000 words into word2Ind and matrix M...
Done.
Running Truncated SVD over 10010 words...
Done.

Question 2.1: Word2Vec Plot Analysis [written] (2 points)

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

```
1 words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum', 'venezuela']
2 plot_embeddings(M_reduced, word2Ind, words)
```



Multiple Choice Question: Why aren't countries "venezuela", "ecuador" and "kuwait" clustered together in the Word2Vec plot while they were clustered together in the co-occurrence plot? - State All That Apply

- A) Word2Vec was trained on a larger dataset in which the countries did not always appear in the same context as the small dataset, used to compute the co-occurrence matrix
- B) The countries are not geographically close to each other

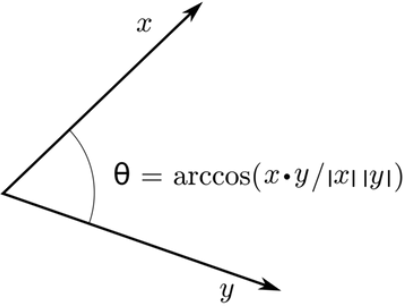
Write your answer here.

A

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: Homophonous Words (2 points) [code + written]

Find a [homophonous](#) 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 homophonous words before you find one.**

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](#).

Written Question: State the homophonous word you discover and the multiple meanings that occur in the top 10. Why do you think many of the homophonous words you tried didn't work?

```
1  ### SOLUTION BEGIN
2
3  wv_from_bin.most_similar("flies")
4
5  ### SOLUTION END
```

```
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
    if np.issubdtype(vec.dtype, np.int):
(['flys', 0.6562365293502808),
 ('fly', 0.6233823895454407),
 ('flying', 0.606692910194397),
 ('flew', 0.5493152141571045),
 ('Pedro_Forteza_pilot', 0.5440078973770142),
 ('Drosophila_melanogaster_fruit', 0.4980913996696472),
 ('aborts_takeoff', 0.4658214747905731),
 ('Flies', 0.46035057306289673),
 ('C._Lubanski', 0.4572356939315796),
 ('puddle_jumper', 0.45273661613464355)]
```

Write your answer here.

The word "flies" worked for me. It gives both "flying" and "Drosophila_melanogaster_fruit". Words such as "rose" and "type" did not work for me. It seems that the two meanings of flies are very closed, but the two meanings of rose/type are very fart apart.

Question 2.3: Synonyms & Antonyms (2 points) [code + written]

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

Find three words (w1,w2,w3) where w1 and w2 are synonyms and w1 and w3 are antonyms, but Cosine Distance(w1,w3) < Cosine Distance(w1,w2). For example, w1="happy" is closer to w3="sad" than to w2="cheerful". (1 point)

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.

Written Question: State which three words (w1, w2, w3) you found. What are some possible explanations for why this counterintuitive result happened?

```
1  ### SOLUTION BEGIN
2
3  w1 = "fast"
4  w2 = "speedy"
5  w3 = "slow"
6  w1_w2_dist = wv_from_bin.distance(w1, w2)
7  w1_w3_dist = wv_from_bin.distance(w1, w3)
8
9  print("Synonyms {}, {} have cosine distance: {}".format(w1, w2, w1_w2_dist))
10 print("Antonyms {}, {} have cosine distance: {}".format(w1, w3, w1_w3_dist))
11
12 ### SOLUTION END
```

```
↳ Synonyms fast, speedy have cosine distance: 0.5282870233058929
Antonyms fast, slow have cosine distance: 0.4686307907104492
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
if np.issubdtype(vec.dtype, np.int):
```

Write your answer here.

The antonyms fast, slow are closer than the synonyms fast, speedy. Possibly because fast and slow appear in the same context more often than fast and speedy.

▼ 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](#).

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

↳ /usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
if np.issubdtype(vec.dtype, np.int):
[('queen', 0.7118192911148071),
 ('monarch', 0.6189674139022827),
 ('princess', 0.5902431011199951),
 ('crown_prince', 0.5499460697174072),
 ('prince', 0.5377321243286133),
 ('kings', 0.5236844420433044),
 ('Queen_Consort', 0.5235945582389832),
 ('queens', 0.518113374710083),
 ('sultan', 0.5098593235015869),
 ('monarchy', 0.5087411999702454)]
```

▼ Question 2.4: Finding Analogies [code + written] (2 Points)

Find an example of analogy that holds according to these vectors (i.e. the intended word is ranked top). **Note:** You may have to try many analogies to find one that works!

Written Question: State the successful analogy you found in the form *x:y :: a:b*. If you believe the analogy is complicated, give a short explanation as to why it holds.

```
1 ### SOLUTION BEGIN
2
3 pprint.pprint(wv_from_bin.most_similar(positive=['cat', 'puppy'], negative=['dog']))
4
5 ### SOLUTION END

↳ /usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
if np.issubdtype(vec.dtype, np.int):
[('kitten', 0.7634989619255066),
 ('puppies', 0.7110899090766907),
 ('pup', 0.6929494142532349),
 ('kittens', 0.6888390779495239),
 ('cats', 0.6796488761901855),
 ('kitties', 0.6261522769927979),
 ('tabby', 0.6248785257339478),
 ('feline', 0.6239446401596069),
 ('beagle', 0.5984721779823303),
 ('tortoiseshell_cat', 0.5960987210273743)]
```

Write your answer here.

dog:puppy :: cat:kitten

▼ Question 2.5: Incorrect Analogy [code + written] (2 point)

Find an example of analogy that does **not** hold according to these vectors.

Written Question: State the intended analogy in the form *x:y :: a:b*, and state the (incorrect) value of *b* according to the word vectors.

```
1 ### SOLUTION BEGIN
2
3 pprint.pprint(wv_from_bin.most_similar(positive=['horse', 'puppy'], negative=['dog']))
```

```
4
5  ### SOLUTION END
```

```
➤ /usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
  if np.issubdtype(vec.dtype, np.int):
[('horses', 0.7203024625778198),
 ('colt', 0.6814397573471069),
 ('mare', 0.6575499773025513),
 ('foal', 0.6572946310043335),
 ('racehorse', 0.6455104947090149),
 ('stallion', 0.6454670429229736),
 ('filly', 0.644342303276062),
 ('gelding', 0.6223742961883545),
 ('foals', 0.5980200171470642),
 ('gelded', 0.592192530632019)]
```

Write your answer here.

The intended analogy is dog:puppy :: horse:mare. Instead word2vec gives horses.

▼ Question 2.6: Guided Analysis of Bias in Word Vectors [written] (1 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 "boss" and most dissimilar to "man", and (b) which terms are most similar to "man" and "boss" and most dissimilar to "woman".

Written Question: Point out one difference between the list of female-associated words and the list of male-associated words, and explain how it is reflecting a potential gender bias. Where the search has returned named entities, you may want to perform a quick search to understand what/who they refer to.

```
1  # Run this cell
2  # Here `positive` indicates the list of words to be similar to and `negative` indicates the list of words to be
3  # most dissimilar from.
4  pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'boss'], negative=['man']))
5  print()
6  pprint.pprint(wv_from_bin.most_similar(positive=['man', 'boss'], negative=['woman']))
```

```
➤ /usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
  if np.issubdtype(vec.dtype, np.int):
[('bosses', 0.5522644519805908),
 ('manageress', 0.49151360988616943),
 ('exec', 0.459408164024353),
 ('Manageress', 0.45598435401916504),
 ('receptionist', 0.4474116861820221),
 ('Jane_Danson', 0.44480547308921814),
 ('Fiz_Jennie_McAlpine', 0.44275766611099243),
 ('Coronation_Street_actress', 0.44275569915771484),
 ('supremo', 0.4409852921962738),
 ('coworker', 0.4398624897003174)]

[('supremo', 0.6097397804260254),
 ('MOTHERWELL_boss', 0.5489562153816223),
 ('CARETAKER_boss', 0.5375303626060486),
 ('YEOVIL_Town_boss', 0.5321705341339111),
 ('head_honcho', 0.5281980037689209),
 ('manager_Stan_Ternent', 0.525971531867981),
 ('Viv_Busby', 0.5256163477897644),
 ('striker_Gabby_Agbonlahor', 0.5250812768936157),
 ('BARNSELEY_boss', 0.5238943099975586),
 ('WIGAN_boss', 0.5175146460533142)]
```

Write your answer here.

The female associated words are much less about "boss", meanwhile the male associated words are mostly about "boss" (football coaches in UK). This reflects a potential gender bias that "woman" is less likely to be correlated with "boss".

▼ Question 2.7: Independent Analysis of Bias in Word Vectors [code + written] (2 points)

Use the most_similar function to find another case where some bias is exhibited by the vectors.

Written Question: Briefly explain the additional example of bias that you discover.

```
1  ### SOLUTION BEGIN
2
3  pprint.pprint(wv_from_bin.most_similar(positive=['chinese','rice'], negative=['american']))
4  print()
5  pprint.pprint(wv_from_bin.most_similar(positive=['american','rice'], negative=['chinese']))
6
```

7 ### SOLUTION END

```
➦ /usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of is
    if np.issubdtype(vec.dtype, np.int):
[('noodles', 0.5879100561141968),
 ('glutinous_rice', 0.58396315574646),
 ('unhusked_rice', 0.5691326856613159),
 ('kangkong', 0.5619485974311829),
 ('fried_noodles', 0.5542783737182617),
 ('rice_porridge', 0.5532837510108948),
 ('Koshihikari', 0.5530127882957458),
 ('unmilled_rice', 0.5502282381057739),
 ('curry_puffs', 0.5494765043258667),
 ('Hom_Mali', 0.5466756224632263)]

[('wheat_flour', 0.5167709589004517),
 ('wheat', 0.5091803669929504),
 ('maize', 0.4864976406097412),
 ('milled_rice', 0.4833839237689972),
 ('corn', 0.4829416275024414),
 ('cotton', 0.48142778873443604),
 ('flour', 0.47442469000816345),
 ('sorghum', 0.46946340799331665),
 ('Basmati_rice', 0.46936652064323425),
 ('lentils', 0.46657830476760864)]
```

Write your answer here.

The Chinese associated words about rice is indeed mostly about different types of rice and some Asian cusines. The American associated words about rice contains many grain products other than rice.

This reflects a potential bias that, rice is perceived as the main source of carbs for Chinese, but not for the Americans.

▼ Question 2.8: Thinking About Bias [written] (1 point)

Multiple Choice Question: What factors might contribute to the biases observed in the word vectors? (State all that apply)

- A) These biases may come from bias in the text (possibly held by people generating the text) that was used as training data
- B) The training corpus may have many instances of sentences that relate certain races, genders, etc. to certain properties or behaviours
- C) These biases come from the dictionary meaning of words

Write your answer here.

A, B, C

Submission Instructions

1. Please make sure you have entered your name and SCPD XID Number above
2. Click the Save button at the top (“File > Save”)
3. Select “Edit > Clear All Outputs”. This will clear all the outputs from all cells (but will keep the content of all cells).
4. Select “Runtime > Run All”. This will run all the cells in order, and will take several minutes.
5. Once you've rerun everything, select “File > Print > Save as PDF”
6. Look at the PDF file and make sure all your solutions are there, displayed correctly including the output cells. The PDF is the only thing your graders will see!
7. Submit your PDF via the Gradescope submission link in the Assignment 1 block of your SCPD learning portal. The system will ask you to take 1-2 minutes to tag pages of your PDF to the corresponding response item. You can see a demonstration of this submission process at this link - <https://youtu.be/yocclo79qh4>

