

CSCI 4140: Natural Language Processing

CSCI/DASC 6040: Computational Analysis of Natural Languages

Spring 2023

Homework 3 - Exploring word vectors

Due Sunday, February 26, at 11:59 PM

Do not redistribute without the instructor's written permission.

```
In [32]: # ALL Import Statements Defined Here
# Note: Do not add to this list.

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

from 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)
# -----
```

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 (40 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 [Word embedding](#)).

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:

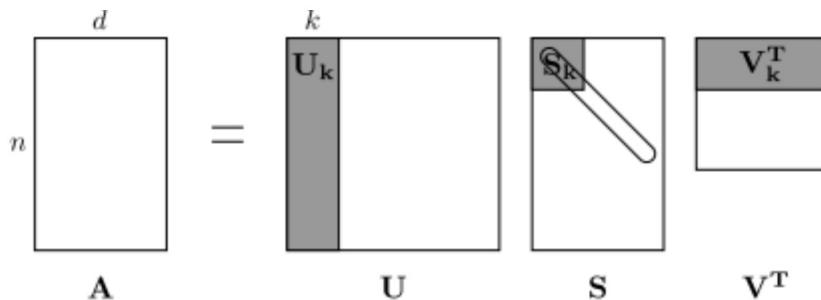
Document 1: "all that glitters is not gold"

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 this case we imagine START and END tokens encapsulating each document, e.g., "START All that glitters is not gold END", and include these tokens in our co-occurrence counts.

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



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

Notes: If you can barely remember what an eigenvalue is, here's [a slow, friendly introduction to SVD](#). Though, for the purpose of this class, you only need to know how to extract the k -dimensional embeddings by utilizing pre-programmed implementations of these algorithms from the numpy, scipy, or sklearn python packages. In practice, it is challenging to apply full SVD to large corpora because of the memory needed to perform PCA or SVD. However, if you only want the top k vector components for relatively small k — known as [Truncated SVD](#) — then there are reasonably scalable techniques to compute those iteratively.

Plotting Co-Occurrence Word Embeddings

Here, we will be using the Reuters (business and financial news) corpus. If you haven't run the import cell at the top of this page, please run it now (click it and press SHIFT-RETURN). The corpus consists of 10,788 news documents totaling 1.3 million words. These documents span 90 categories and are split into train and test. For more details, please see <https://www.nltk.org/book/ch02.html>. We provide a `read_corpus` function below that pulls out only articles from the "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.

```
In [33]: 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....

```
In [34]: reuters_corpus = read_corpus()
pprint.pprint(reuters_corpus[:3], compact=True, width=100)
```

[['<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand', 'dow
nwards', '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', 'projec
tion', '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', 'jap
an', "", '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', '(', '2
1', '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', 'stre
amlined', 'u',
 '.', 's', '.', 'petrochemical', 'industry', 'to', 'record', 'profits', 'this', 'y
ear', ',',
 'with', 'growth', 'expected', 'through', 'at', 'least', '1990', ',', 'major', 'co
mpany',
 'executives', 'predicted', '.', 'this', 'bullish', 'outlook', 'for', 'chemical',
 'manufacturing',
 'and', 'an', 'industrywide', 'move', 'to', 'shed', 'unrelated', 'businesses', 'ha
s', 'prompted',
 'gaf', 'corp', '&', 'lt', ';', 'gaf', '>', 'privately', '-', 'held', 'cain', 'ch
emical', 'inc',
 ', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions', 'of',
 'petrochemical',
 'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&', 'l
t', ';', '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'

lly', 'cheaper',
 'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was', 'weakenin
g', 'against',
 'foreign', 'currencies', '.', 'that', 'helped', 'boost', 'u', '.', 's', '.', 'che
mical',
 'exports', '.', 'also', 'helping', 'to', 'bring', 'supply', 'and', 'demand', 'int
o', 'balance',
 'has', 'been', 'the', 'gradual', 'market', 'absorption', 'of', 'the', 'extra', 'c
hemical',
 'manufacturing', 'capacity', 'created', 'by', 'middle', 'eastern', 'oil', 'produc
ers', 'in',
 'the', 'early', '1980s', '.', 'finally', ',', 'virtually', 'all', 'major', 'u',
'.', 's', '.',
 'chemical', 'manufacturers', 'have', 'embarked', 'on', 'an', 'extensive', 'corpor
ate',
 'restructuring', 'program', 'to', 'mothball', 'inefficient', 'plants', ',', 'trim
, 'the',
 'payroll', 'and', 'eliminate', 'unrelated', 'businesses', '.', 'the', 'restructur
ing', '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', 'b
illion', 'dlrs',
 'for', 'borg', 'warner', 'corp', '&', 'lt', ';', 'bor', '>', 'a', 'chicago', 'ma
nufacturer',
 'of', 'plastics', 'and', 'chemicals', '.', 'another', 'industry', 'powerhouse',
, ', 'w', '.',
 'r', '.', 'grace', '&', 'lt', ';', 'gra', '>', 'has', 'divested', 'its', 'retaili
ng', ',',
 'restaurant', 'and', 'fertilizer', 'businesses', 'to', 'raise', 'cash', 'for', 'c
hemical',
 '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', 'commoditi
es', ',', 'such',
 'as', 'ethylene', ',', 'in', 'favor', 'of', 'more', 'profitable', 'specialty', 'c
hemicals',
 'that', 'are', 'custom', '-', 'designed', 'for', 'a', 'small', 'group', 'of', 'bu
yers', '.', "",
 'companies', 'like', 'dupont', '&', 'lt', ';', 'dd', '>', 'and', 'monsanto', 'co
, '&', 'lt', ';',
 'mtc', '>', 'spent', 'the', 'past', 'two', 'on', 'three', 'years', 'trying', 'to
, 'get', 'out',
 'of', 'the', 'commodity', 'chemical', 'business', 'in', 'reaction', 'to', 'how',
'badly', 'the',
 'market', 'had', 'deteriorated', ',"', 'dosher', 'said', '.', "", 'but', 'i', 't
hink', 'they',
 'will', 'eventually', 'kill', 'the', 'margins', 'on', 'the', 'profitable', 'chemi
cals', 'in',
 'the', 'niche', 'market', '."', 'some', 'top', 'chemical', 'executives', 'share',
'the',
 'concern', '.', "", 'the', 'challenge', 'for', 'our', 'industry', 'is', 'to', 'k
eep', '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', '.', "", 'a',
 'the',
 'issue', 'is', 'basicly', '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
>']]
```

Question 1.1: Implement distinct_words [code] (8 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.

```
In [35]: import itertools
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
        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(itertools.chain(*corpus))))
    num_corpus_words = len(corpus_words)

    return corpus_words, num_corpus_words
```

```
In [36]: # -----
# 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 we  
test_corpus_words, num_corpus_words = distinct_words(test_corpus)  
  
# Correct answers
ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that", "gold", "A  
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 wor  
  
# Test correct words
assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.\nCorr  
  
# Print Success
print ("-" * 80)
print("Passed All Tests!")
print ("-" * 80)  
-----  
Passed All Tests!
```

Question 1.2: Implement `compute_co_occurrence_matrix` [code] (12 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 [Python NumPy tutorial](#).

```
In [37]: from collections import Counter

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

    Note: Each word in a document should be at the center of a window. Words need not be adjacent.
    number of co-occurring words.

    For example, if we take the document "START All that glitters is not All"
    "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
    Returns:
        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
            word2Ind (dict): dictionary that maps word to index (i.e. row/column number)
    """
    words, num_words = distinct_words(corpus)
    M = None
    word2Ind = {}

    # -----
    # Write your implementation here.
    # -----
    M = np.zeros((num_words, num_words))

    coocc = {i:Counter({j:0 for j in words if j!=i}) for i in words} #

    index = 0
    for key in coocc.keys():
        word2Ind.update({key : index})
        index += 1

    for sen in corpus:
        for i in range(len(sen)):
            if i < window_size:
                c = Counter(sen[0:i+window_size+1])
                del c[sen[i]]
                coocc[sen[i]] = coocc[sen[i]] + c
            elif i > len(sen)-(window_size+1):
                c = Counter(sen[i-window_size::])
                del c[sen[i]]
                coocc[sen[i]] = coocc[sen[i]] + c
            else:
                c = Counter(sen[i-window_size:i+window_size+1])
                del c[sen[i]]
                coocc[sen[i]] = coocc[sen[i]] + c

    for i in coocc.keys():
        for j in coocc[i].keys():
            M[word2Ind[i],word2Ind[j]]+=coocc[i][j]
```

```
    return M, word2Ind
```

```
In [38]: # -----
# 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 we
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., 1., 0.,],
     [0., 0., 1., 0., 0., 0., 0., 0., 1.,],
     [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., 1., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 1., 1.,],
     [0., 0., 1., 0., 0., 0., 1., 0., 0.,],
     [0., 0., 0., 0., 1., 1., 0., 0., 0.,],
     [1., 0., 0., 1., 1., 0., 0., 0., 1.,],
     [0., 1., 1., 0., 1., 0., 0., 1., 0.,]])
)
word2Ind_ans = {'All': 0, "All's": 1, 'END': 2, 'START': 3, 'ends': 4, 'glitters':  
  
# Test correct word2Ind
assert (word2Ind_ans == word2Ind_test), "Your word2Ind is incorrect:\nCorrect: {}\n  
  
# Test correct M shape
assert (M_test.shape == M_test_ans.shape), "M matrix has incorrect shape.\nCorrect:  
  
# 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 mat  
  
# Print Success
print("-" * 80)
print("Passed All Tests!")
print("-" * 80)
```

Passed All Tests!

Question 1.3: Implement `reduce_to_k_dim` [code] (4 points)

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`.

```
In [39]: def reduce_to_k_dim(M, k=2):
    """ Reduce a co-occurrence count matrix of dimensionality (num_corpus_words, num
        to a matrix of dimensionality (num_corpus_words, k) using the following SVD
        - http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.Truncat
    """
    Params:
        M (numpy matrix of shape (number of corpus words, number of corpus word
        k (int): embedding size of each word after dimension reduction
    Return:
        M_reduced (numpy matrix of shape (number of corpus words, k)): matrix o
            In terms of the SVD from math class, this actually returns U *
    """
    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
```

```
In [40]: # -----
# 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 we
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])
assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have {}".format(M_test_reduced.shape[1])

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

Running Truncated SVD over 10 words...

Done.

Passed All Tests!

Question 1.4: Implement plot_embeddings [code] (4 points)

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.

```
In [41]: def plot_embeddings(M_reduced, word2Ind, words):
    """ Plot in a scatterplot the embeddings of the words specified in the list "wo
    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
        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:
        x = M_reduced[word2Ind[word]][0]
        y = M_reduced[word2Ind[word]][1]
        plt.scatter(x, y, color='red')
        plt.text(x, y, word, fontsize=9)

    plt.show()
```

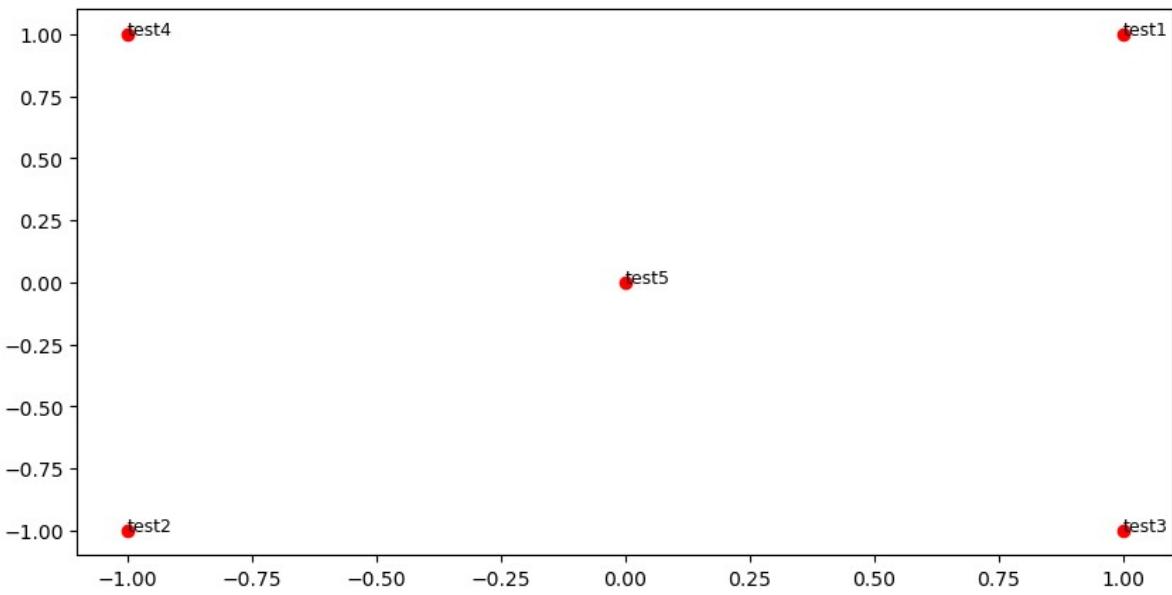
```
In [42]: # -----
# Run this sanity check
# Note that this NOT an exhaustive check for correctness.
# The plot produced should look like the "test solution plot" depicted below.
# -----

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

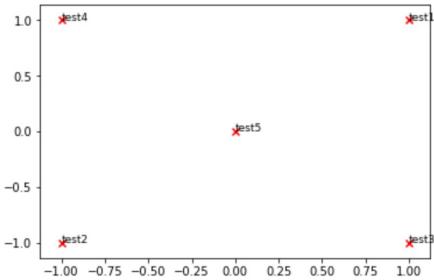
M_reduced_plot_test = np.array([[1, 1], [-1, -1], [1, -1], [-1, 1], [0, 0]])
word2Ind_plot_test = {'test1': 0, 'test2': 1, 'test3': 2, 'test4': 3, 'test5': 4}
words = ['test1', 'test2', 'test3', 'test4', 'test5']
plot_embeddings(M_reduced_plot_test, word2Ind_plot_test, words)

print ("-" * 80)
```

Outputted Plot:



****Test Plot Solution****



Question 1.5: Co-Occurrence Plot Analysis [written] (12 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 $U \times S$, 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.

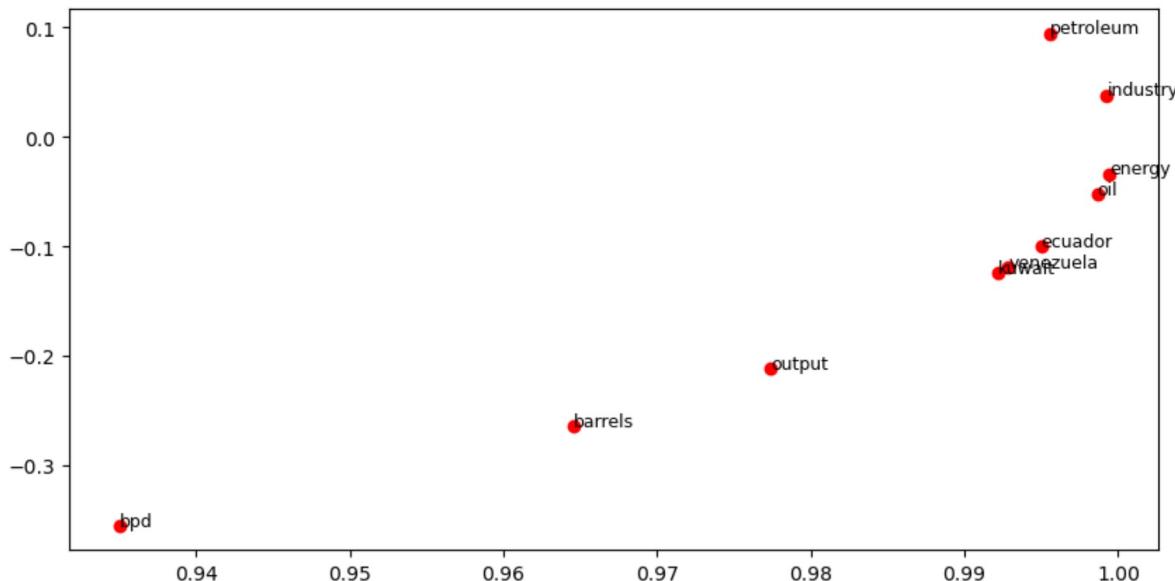
```
In [43]: # -----
# Run This Cell to Produce Your Plot
# -----
reuters_corpus = read_corpus()
M_co_occurrence, word2Ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

# 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']
plot_embeddings(M_normalized, word2Ind_co_occurrence, words)
```

Running Truncated SVD over 8185 words...

Done.



Kuwait, Venezuela, and Ecuador all cluster together. This makes sense to me because they are all countries and may appear in a list together. Words like bpd, output, and industry didn't cluster together. This seems counter to me because they seem like they would be used together, i.e. "a 100,000 bpd output".

Part 2: Prediction-Based Word Vectors (60 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 might take several minutes.

```
In [44]: 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.key_to_index.keys())
    print("Loaded vocab size %i" % len(vocab))
    return wv_from_bin
```

```
In [45]: # -----
# Run Cell to Load Word Vectors
# Note: This may take several minutes
# -----
wv_from_bin = load_word2vec()
```

Loaded vocab size 3000000

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, and
2. Reduce the vectors from 300-dimensional to 2-dimensional by running `reduce_to_k_dim` (your Truncated SVD function).

```
In [46]: def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd', 'ecuador',
    """ Put the word2vec vectors into a matrix M.
    Param:
        wv_from_bin: KeyedVectors object; the 3 million word2vec vectors loaded
    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.key_to_index.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.get_vector(w))
            word2Ind[w] = curInd
            curInd += 1
        except KeyError:
            continue
    for w in required_words:
        try:
            M.append(wv_from_bin.get_vector(w))
            word2Ind[w] = curInd
            curInd += 1
        except KeyError:
            continue
    M = np.stack(M)
    print("Done.")
    return M, word2Ind
```

```
In [47]: # -----
# Run Cell to Reduce 300-Dimensinal 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)
```

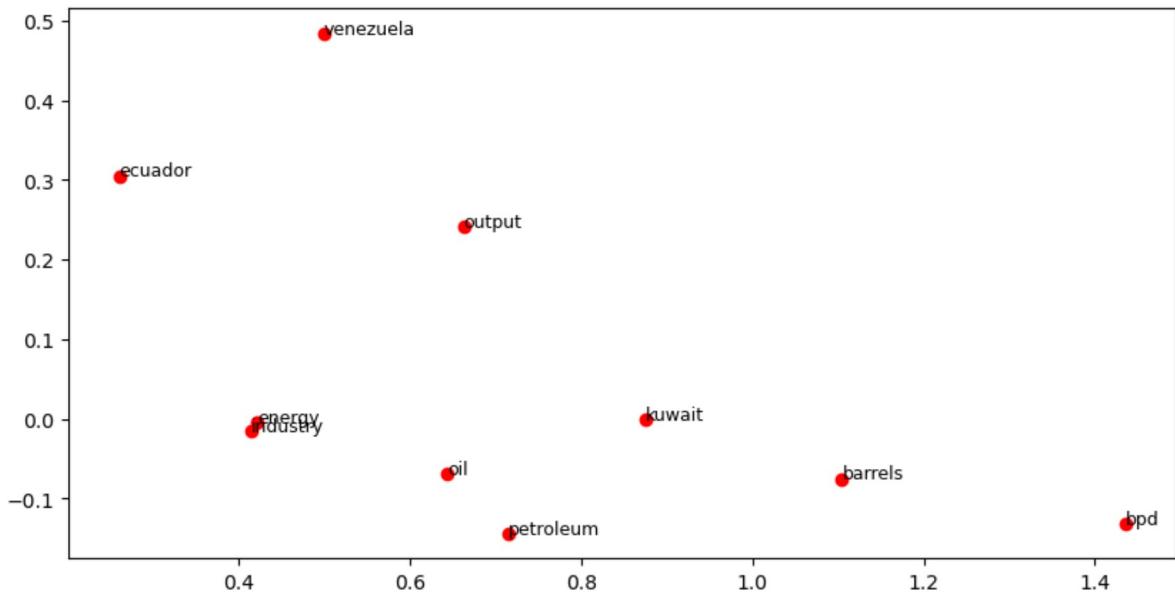
Shuffling words ...
Putting 10000 words into word2Ind and matrix M...
Done.
Running Truncated SVD over 10010 words...
Done.

Question 2.1: Word2Vec Plot Analysis [written] (16 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?

```
In [48]: words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output',  
plot_embeddings(M_reduced, word2Ind, words)
```

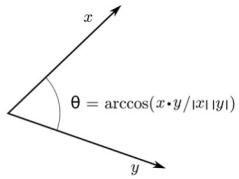


The words industry and energy are very close, while oil and petroleum are fairly close. The countries didn't cluster together like I thought they would. This plot seems to be more spread out and doesn't have the same shape as the one before.

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: Polysemous Words (8 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](#).

```
In [49]: # -----
# Write your polysemous word exploration code here.
temp = wv_from_bin.most_similar("spring")
print("wv_from_bin.most_similar(\"spring\"):")
for thing in temp:
    print(thing[0])

temp = wv_from_bin.most_similar("base")
print('\n')
print("wv_from_bin.most_similar(\"base\"):")
for thing in temp:
    print(thing[0])

temp = wv_from_bin.most_similar("bark")
print('\n')
print("wv_from_bin.most_similar(\"bark\"):")
for thing in temp:
    print(thing[0])

# -----
```

```
wv_from_bin.most_similar("spring"):  
summer  
winter  
autumn  
springtime  
midsummer  
Punxsutawney_Phil_predicts  
Spring  
week  
fall  
summertime
```

```
wv_from_bin.most_similar("base"):  
bases  
roller_skates_whirled  
Base  
Medidata_diverse  
Hickam_Air_Force  
Mehran_naval_aviation  
Pierzynski_hustled  
QUIOCHO_singled  
Camp_Ederle  
Soffront_installed
```

```
wv_from_bin.most_similar("bark"):  
barks  
Pacific_yew_tree  
barking  
cork_oak_tree  
beetles_burrow  
barky  
cambium  
sapwood  
frass  
barked
```

Bark is a polysemous word, it contains barked the past tense of bark the action and barky meaning covered with or resembling bark a reference to wood. I think that many of the polysemous words that I tried didn't work because many of the associated words were derivations of the word but didn't change the meaning of the word.

Question 2.3: Synonyms & Antonyms (8 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 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.

```
In [50]: # -----
# Write your synonym & antonym exploration code here.

#w1 = "happy"
#w2 = "cheerful"
#w3 = "sad"
w1 = "clean"
w2 = "fresh"
w3 = "dirty"
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 clean, fresh have cosine distance: 0.6760035157203674
Antonyms clean, dirty have cosine distance: 0.4819817543029785

I think that one possibility for this example is that the words clean and dirty probably appear "together" more often than clean and fresh. This might be because fresh is another word for clean and therefore will not appear "together" with the word clean. However this isn't true for dirty, it would likely appear with clean, such as "how to clean that dirty oven."

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](#).

In [51]:

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

[('queen', 0.7118193507194519),
 ('monarch', 0.6189674735069275),
 ('princess', 0.5902431011199951),
 ('crown_prince', 0.5499460697174072),
 ('prince', 0.5377322435379028),
 ('kings', 0.5236844420433044),
 ('Queen_Consort', 0.5235945582389832),
 ('queens', 0.5181134939193726),
 ('sultan', 0.5098593235015869),
 ('monarchy', 0.5087411403656006)]
```

Question 2.4: Finding Analogies [code + written] (8 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!

In [52]:

```
# -----
# Write your analogy exploration code here.

pprint.pprint(wv_from_bin.most_similar(positive=['cat', 'puppy'], negative=['dog']))

# -----
```

```
[('kitten', 0.7634989619255066),  
 ('puppies', 0.7110899090766907),  
 ('pup', 0.6929495334625244),  
 ('kittens', 0.6888389587402344),  
 ('cats', 0.6796489357948303),  
 ('kitties', 0.6261522769927979),  
 ('tabby', 0.6248785257339478),  
 ('feline', 0.6239445805549622),  
 ('beagle', 0.5984722375869751),  
 ('tortoiseshell_cat', 0.5960987210273743)]
```

dog:puppy::cat:kitten

Question 2.5: Incorrect Analogy [code + written] (4 points)

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.

```
In [53]: # -----  
# Write your incorrect analogy exploration code here.  
pos = ['leopard', 'dog']  
neg = ['wolf']  
pprint.pprint(wv_from_bin.most_similar(positive=pos, negative=neg))  
  
# -----  
  
[('dogs', 0.6001646518707275),  
 ('schnauzer', 0.59073406457901),  
 ('Shih_Tzu', 0.5869850516319275),  
 ('puppy', 0.5866811275482178),  
 ('Golden_Retriever', 0.5854203701019287),  
 ('Pomeranian', 0.584109902381897),  
 ('bulldog', 0.581149160861969),  
 ('Doberman', 0.5709829926490784),  
 ('shih_tzu', 0.5677027106285095),  
 ('pit_bull_mix', 0.5656912922859192)]
```

wolf:dog :: leopard:cat. Instead we got dogs.

Question 2.6: Guided Analysis of Bias in Word Vectors [written] (4 points)

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". What do you find in the top 10?

```
In [54]: # Run this cell
# Here `positive` indicates the list of words to be similar to and `negative` indicates the most dissimilar from.
pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'boss'], negative=['man']))
print()
pprint.pprint(wv_from_bin.most_similar(positive=['man', 'boss'], negative=['woman']))

[('bosses', 0.5522643327713013),
 ('manageress', 0.49151360988616943),
 ('exec', 0.45940810441970825),
 ('Manageress', 0.4559844434261322),
 ('receptionist', 0.4474117159843445),
 ('Jane_Danson', 0.44480544328689575),
 ('Fiz_Jennie_McAlpine', 0.4427577257156372),
 ('Coronation_Street_actress', 0.44275563955307007),
 ('supremo', 0.4409853518009186),
 ('coworker', 0.43986251950263977)]

[('supremo', 0.6097397804260254),
 ('MOTHERWELL_boss', 0.5489562153816223),
 ('CARETAKER_boss', 0.5375303030014038),
 ('Bully_Wee_boss', 0.5333974957466125),
 ('YEOVIL_Town_boss', 0.5321704745292664),
 ('head_honcho', 0.5281979441642761),
 ('manager_Stan_Ternent', 0.525971531867981),
 ('Viv_Busby', 0.5256163477897644),
 ('striker_Gabby_Agbonlahor', 0.5250813961029053),
 ('BARNESLEY_boss', 0.5238943696022034)]
```

The words that are most similar to woman and boss but are dissimilar to man are: bosses, manageress, exec, and receptionist. The words that are most similar to man and boss but are dissimilar to woman are: supremo, motherwell_boss, caretaker_boss, Bully_wee_boss, YEOVIL_Town_boss, and head_honcho. While some of these tokens I'm just guessing at, it seems like the words that are similar to woman-boss, and not man are not really considering the woman the boss. Receptionist makes the most clear cut example. Where as many of the tokens for the man boss combination seem to be examples of men who are bosses, or another word that would mean boss.

Question 2.7: Independent Analysis of Bias in Word Vectors [code + written] (8 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.

```
In [55]: # -----
# Write your bias exploration code here.

pprint.pprint(wv_from_bin.most_similar(positive=['chef', 'man'], negative=['woman']))
print()
pprint.pprint(wv_from_bin.most_similar(positive=['chef', 'woman'], negative=['man']))

# -----
```

```
[('Chef', 0.7118616104125977),
 ('sous_chef', 0.6578055024147034),
 ('Executive_Chef', 0.6405356526374817),
 ('chefs', 0.6211988925933838),
 ('pastry_chef', 0.6104018688201904),
 ('Michelin_starred_chef', 0.5813212394714355),
 ('restaurateur', 0.571975827217102),
 ('pizzaiolo', 0.5600727796554565),
 ('celebrity_chef', 0.5592496395111084),
 ('chef_proprietor', 0.5564097166061401)]
```



```
[('pastry_chef', 0.7035524249076843),
 ('chefs', 0.667602002620697),
 ('sous_chef', 0.651990532875061),
 ('Chef', 0.646479606628418),
 ('Pastry_Chef', 0.5879147052764893),
 ('celebrity_chef', 0.5874067544937134),
 ('Suvir_Saran', 0.5783470869064331),
 ('George_Mavrothalassitis', 0.5763151049613953),
 ('Michelin_starred_chef', 0.5738801956176758),
 ('Donatella_Arpaia', 0.5737050175666809)]
```

When searching through different possibilities I noticed that if one of the positive words is something that can be "easily" divided into groups based on some binary category it was easier to find some bias in the results. Chef returns terms like restaurateur, and chef_proprietor with man associated, indicating that a man would own the restaurant he works in. While woman's first item is pastry_chef, and the list is missing the other positive terms from the male category.

Question 2.8: Thinking About Bias [written] (4 points)

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

As indicated above, it seems that words that have easily definable binary associations, are more commonly associated with bias. This would show up in the word vectors since these words are probably used in context with their respected binary category in the content that we are building the model on. Since these articles are written by people from a place at a certain time from a certain culture their biases will likely be incorporated into the word combinations they form when speaking about a certain topic. For instance a woman writing an article about ballerinas may not address men at all, even though men do dance in ballet performance.

How to submit this problem set:

- Write all the answers in this iPython notebook. Once you are finished (1) generate the PDF file (`File -> Print Preview`, and print to PDF), 2) ZIP the PDF and this Jupyter Notebook (`.ipynb`), and 3) upload the ZIP file to Canvas.
- **Important:** check your PDF before you turn it in to Canvas to make sure it exported correctly.
- When creating your final version of the PDF to hand in, please do a fresh restart and execute every cell in order. Then you'll be sure it's actually right. One handy way to do this is by clicking `Runtime -> Run All` in the notebook menu.

In []:

In []:

In []:

In []: