

# exploring\_word\_vectors

January 10, 2019

## 1 CS224N Assignment 1: Exploring Word Vectors (25 Points)

Welcome to CS224n!

Before you start, make sure you read the README.txt in the same directory as this notebook.

```
In [1]: # All Import Statements Defined Here
        # Note: Do not add to this list.
        # All the dependencies you need, can be installed by running .
        # -----
```

```
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)
# -----
```

```
c:\users\agoswami\appdata\local\continuum\anaconda3\envs\py36\lib\site-packages\gensim\utils.p
warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

```
[nltk_data] Downloading package reuters to
[nltk_data] C:\Users\agoswami\AppData\Roaming\nltk_data...
[nltk_data] Package reuters is already up-to-date!
```

## 1.1 Word Vectors

Word Vectors are often used as a fundamental component for downstream NLP tasks, e.g. question answering, text generation, translation, etc., so it is important to build some intuitions as to their strengths and weaknesses. Here, you will explore two types of word vectors: those derived from *co-occurrence matrices*, and those derived via *word2vec*.

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

**Note on Terminology:** The terms "word vectors" and "word embeddings" are often used interchangeably. The term "embedding" refers to the fact that we are encoding aspects of a word's meaning in a lower dimensional space. As [Wikipedia](#) states, "*conceptually it involves a mathematical embedding from a space with one dimension per word to a continuous vector space with a much lower dimension*".

## 1.2 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](#)).

### 1.2.1 Co-Occurrence

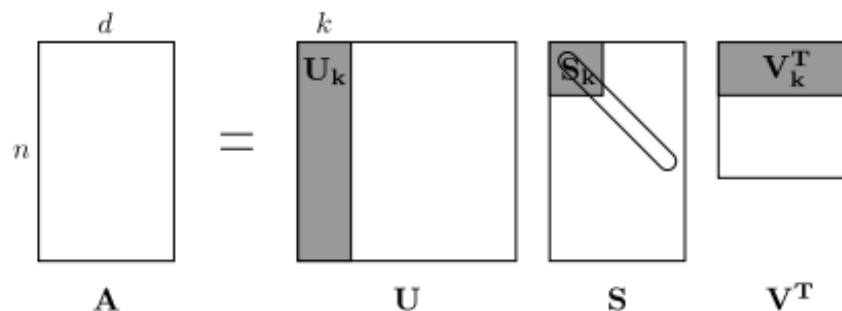
A co-occurrence matrix counts how often things co-occur in some environment. Given some word  $w_i$  occurring in the document, we consider the *context window* surrounding  $w_i$ . Supposing our fixed window size is  $n$ , then this is the  $n$  preceding and  $n$  subsequent words in that document,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

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

Document 1: "all that glitters is not gold"

Document 2: "all is well that ends well"

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



Picture of an SVD

*	START	all	that	glitters	is	not	gold	well	ends	END
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](#). 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.

### 1.2.2 Plotting Co-Occurrence Word Embeddings

Here, we will be using the Reuters (business and financial news) corpus. If you haven't run the import cell at the top of this page, please run it now (click it and press SHIFT-

RETURN). The corpus consists of 10,788 news documents totaling 1.3 million words. These documents span 90 categories and are split into train and test. For more details, please see <https://www.nltk.org/book/ch02.html>. We provide a `read_corpus` function below that pulls out only articles from the "crude" (i.e. news articles about oil, gas, etc.) category. The function also adds START and END tokens to each of the documents, and lowercases words. You do **not** have to perform any other kind of pre-processing.

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

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

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

```
[['<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand', 'downwards', 'the',
  'ministry', 'of', 'international', 'trade', 'and', 'industry', '(', 'miti', ')', 'will', 'review',
  'its', 'long', '-', 'term', 'energy', 'supply', '/', 'demand', 'outlook', 'by', 'august', 'to',
  'meet', 'a', 'forecast', 'downtrend', 'in', 'japanese', 'energy', 'demand', ',', 'the', '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', 'review',
  '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', '"', 'the',
  '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', 'in
```

', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions', 'of', 'petrochem  
 'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&', 'lt', ';', 'as  
 '>', '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', ',', '"', 's',  
 'paul', 'oreffice', ',', 'chairman', 'of', 'giant', 'dow', 'chemical', 'co', '&', 'lt', ';',  
 'dow', '>', 'adding', ',', '"', 'there', '"', 's', 'no', 'major', 'plant', 'capacity', 'bein  
 'added', 'around', 'the', 'world', 'now', '.', 'the', 'whole', 'game', 'is', 'bringing', 'ou  
 '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', 'pri  
 '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', 'eas  
 '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', 'pc  
 'gain', 'in', 'profits', 'during', '1987', '.', 'last', 'year', ',', 'the', 'domestic',  
 'industry', 'earned', 'a', 'total', 'of', '13', 'billion', 'dlrs', ',', 'a', '54', 'pct', 'l  
 'from', '1985', '.', 'the', 'turn', 'in', 'the', 'fortunes', 'of', 'the', 'once', '-', 'sick  
 'chemical', 'industry', 'has', 'been', 'brought', 'about', 'by', 'a', 'combination', 'of', '  
 'and', 'planning', ',', 'said', 'pace', '"', 's', 'john', 'dosher', '.', 'dosher', 'said', '  
 'year', '"', 's', 'fall', 'in', 'oil', 'prices', 'made', 'feedstocks', 'dramatically', 'cheap  
 'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was', 'weakening', 'again  
 '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', 'touch  
 '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', 'd  
 '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', ',', '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', 'on' 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', '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', 'of', 'the', 'economy', 'and', 'other', 'rights', 'of', 'both', 'countries', '.', '"', 'as', 'the', 'issue', 'is', 'basically', 'political', ',', 'a', 'solution', 'can', 'only', 'be', 'found', 'through' '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', 'an' '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', 'betw
the', 'two', 'countries', ', ', 'the', 'foreign', 'ministry', 'statement', 'said', '.', 'tu
"', 's', 'ambassador', 'in', 'athens', ', ', 'nazmi', 'akiman', ', ', 'was', 'due', 'to', 'me
'prime', 'minister', 'andreas', 'papandreou', 'today', 'for', 'the', 'greek', 'reply', 'to'
'message', 'sent', 'last', 'week', 'by', 'turkish', 'prime', 'minister', 'turgut', 'ozal', '
the', 'contents', 'of', 'the', 'message', 'were', 'not', 'disclosed', '.', '<END>']]
```

### 1.2.3 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.

```
In [4]: def distinct_words(corpus):
        """ Determine a list of distinct words for the corpus.
        Params:
            corpus (list of list of strings): corpus of documents
        Return:
            corpus_words (list of strings): list of distinct words across the corpus,
            num_corpus_words (integer): number of distinct words across the corpus
        """
        corpus_words = []
        num_corpus_words = -1

        # -----
        # Write your implementation here.

        flattened_list = [word for i in range(len(corpus)) for word in corpus[i]]

        corpus_words = sorted(set(flattened_list))
        num_corpus_words = len(corpus_words)

        # -----

        return corpus_words, num_corpus_words

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

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





```

word2Ind = {}

# -----
# Write your implementation here.

for idx, word in enumerate(words):
    word2Ind[word] = idx

M = np.zeros((num_words, num_words))

for line in corpus:
    for i in range(len(line)):
        i_idx = word2Ind[line[i]]

        for j in range(i-window_size, i+window_size+1):

            if j < 0 or j == i or j >= len(line):
                continue

            j_idx = word2Ind[line[j]]
            M[i_idx, j_idx] += 1

# -----

return M, word2Ind

```

```

In [7]: # -----
# Run this sanity check
# Note that this is not an exhaustive check for correctness.
# -----

```

```

# Define toy corpus and get student's co-occurrence matrix

```

```

test_corpus = ["START All that glitters isn't gold END".split(" "), "START All's well t
M_test, word2Ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)

```

```

# Correct M and word2Ind

```

```

M_test_ans = np.array(
    [[0., 0., 0., 1., 0., 0., 0., 0., 0., 1., 0.],
     [0., 0., 0., 1., 0., 0., 0., 0., 0., 0., 1.],
     [0., 0., 0., 0., 0., 0., 1., 0., 0., 0., 1.],
     [1., 1., 0., 0., 0., 0., 0., 0., 0., 0., 0.],
     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1., 0.],
     [0., 0., 0., 0., 0., 0., 0., 1., 1., 0., 0.],
     [0., 0., 1., 0., 0., 0., 0., 1., 0., 0., 0.],
     [0., 0., 0., 0., 0., 1., 1., 0., 0., 0., 0.],
     [1., 0., 0., 0., 1., 1., 0., 0., 0., 0., 1.],
     [0., 1., 1., 0., 1., 0., 0., 0., 1., 0., 0.]]

```

```

)
word2Ind_ans = {'All': 0, "All's": 1, 'END': 2, 'START': 3, 'ends': 4, 'glitters': 5,

# Test correct word2Ind
assert (word2Ind_ans == word2Ind_test), "Your word2Ind is incorrect:\nCorrect: {}\nYour

# 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 matrix

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

```

```

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

```

### 1.2.5 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).

```

In [8]: def reduce_to_k_dim(M, k=2):
        """ Reduce a co-occurrence count matrix of dimensionality (num_corpus_words, num_corpus_words)
            to a matrix of dimensionality (num_corpus_words, k) using the following SVD function from sklearn:
            - http://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD

```

```

Params:
    M (numpy matrix of shape (number of unique words in the corpus , number of
    k (int): embedding size of each word after dimension reduction
Return:
    M_reduced (numpy matrix of shape (number of corpus words, k)): matrix of k
    In terms of the SVD from math class, this actually returns  $U * S$ 
"""
n_iters = 10      # Use this parameter in your call to `TruncatedSVD`
M_reduced = None
print("Running Truncated SVD over %i words..." % (M.shape[0]))

# -----
# Write your implementation here.

svd = TruncatedSVD(n_components=k, n_iter=n_iters)
M_reduced = svd.fit_transform(M)

# -----

print("Done.")
return M_reduced

```

```

In [9]: # -----
# Run this sanity check
# Note that this not an exhaustive check for correctness
# In fact we only check that your M_reduced has the right dimensions.
# -----

# Define toy corpus and run student code
test_corpus = ["START All that glitters isn't gold END".split(" "), "START All's well +
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
assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have {}".form

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

```

Running Truncated SVD over 10 words...  
Done.

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

### 1.2.6 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.

```
In [10]: def plot_embeddings(M_reduced, word2Ind, words):
        """ Plot in a scatterplot the embeddings of the words specified in the list "words".
            NOTE: do not plot all the words listed in M_reduced / word2Ind.
            Include a label next to each point.

            Params:
                M_reduced (numpy matrix of shape (number of unique words in the corpus , dimensionality))
                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 i, word in enumerate(words):

            word_idx = word2Ind[word]
            word_vec = M_reduced[word_idx]

            x = word_vec[0]
            y = word_vec[1]
            plt.scatter(x, y, marker='x', color='red')
            plt.text(x+0.001, y+0.001, word, fontsize=9)

        plt.show()

        # -----

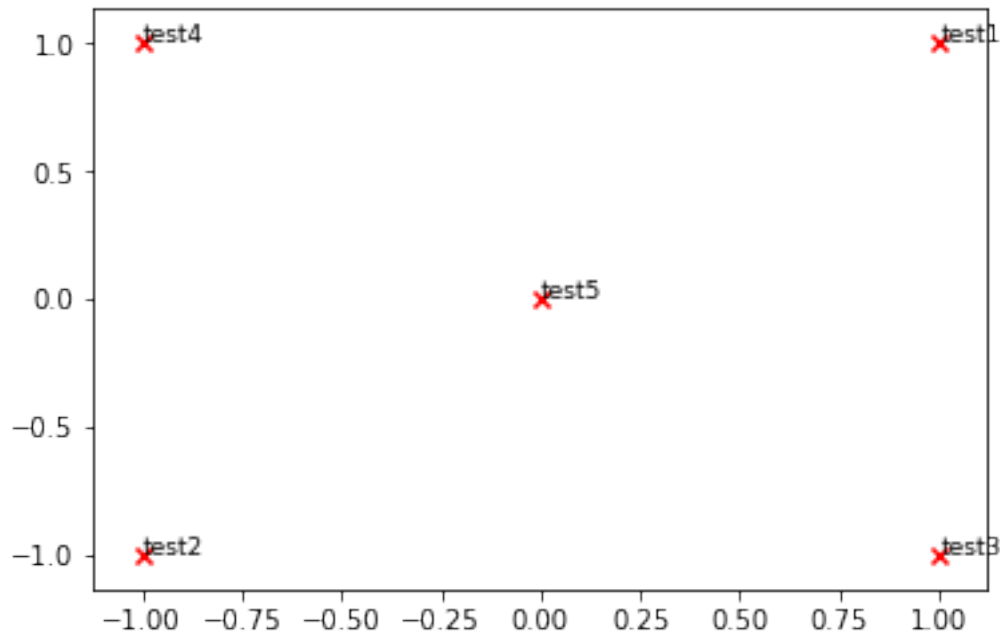
In [11]: # -----
        # Run this sanity check
        # Note that this not an exhaustive check for correctness.
        # The plot produced should look like the "test solution plot" depicted below.
        # -----

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

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

```
print ("-" * 80)
```

Outputted Plot:



### Test Plot Solution

#### 1.2.7 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  $U \cdot S$ , so we normalize the returned vectors, so that all the vectors will appear around the unit circle (therefore closeness is directional closeness). **Note:** The line of code below that does the normalizing uses the NumPy concept of *broadcasting*. If you don't know about broadcasting, check out [Computation on Arrays: Broadcasting by Jake VanderPlas](#).

Run the below cell to produce the plot. It'll probably take a few seconds to run. What clusters together in 2-dimensional embedding space? What doesn't cluster together that you might think should have? **Note:** "bpd" stands for "barrels per day" and is a commonly used abbreviation in crude oil topic articles.

```

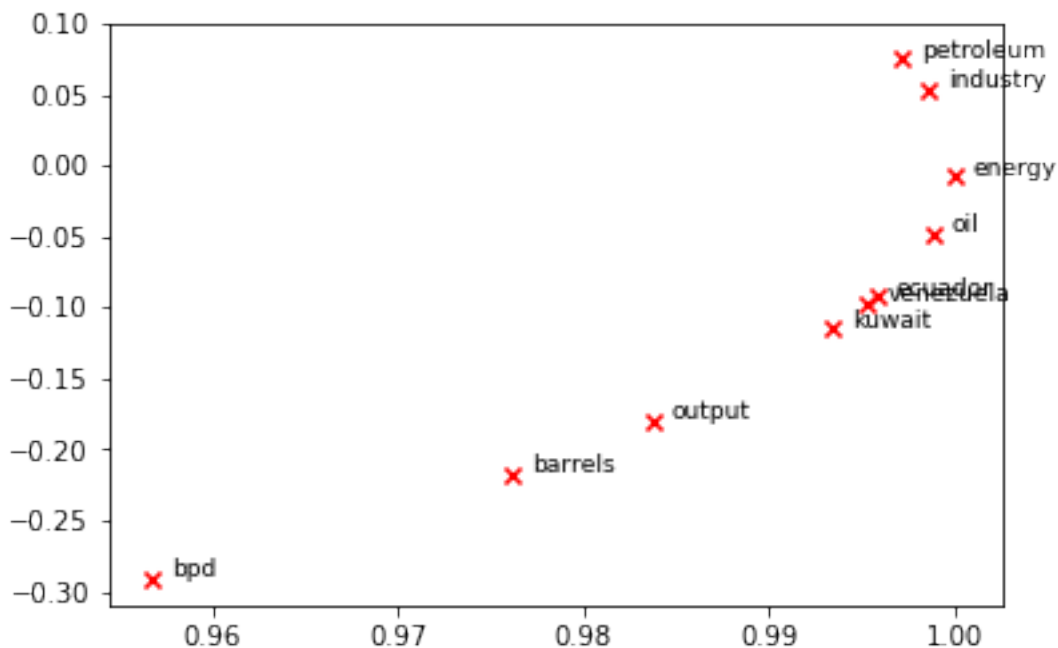
In [12]: # -----
# Run This Cell to Produce Your Plot
# -----
reuters_corpus = read_corpus()
M_co_occurrence, word2Ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

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

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

```

Running Truncated SVD over 8185 words...  
Done.



**Write your answer here.** We observe some clusters in the above plot:

- {petroleum, industry} are terms usually used together. We see them being clustered together. Same with {oil, energy}
- {kuwait, ecuador, venezuela} are the oil producing countries. They are members of the OPEC, and are used frequently together.
- {bpd, barrels, output} are industrial terms, again frequently used together.

Overall, it does seem the clustering is being done based on co-occurrence

### 1.3 Part 2: Prediction-Based Word Vectors (15 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 [13]: def load_word2vec():
         """ Load Word2Vec Vectors
         Return:
             wv_from_bin: All 3 million embeddings, each length 300
         """
         import gensim.downloader as api
         wv_from_bin = api.load("word2vec-google-news-300")
         vocab = list(wv_from_bin.vocab.keys())
         print("Loaded vocab size %i" % len(vocab))
         return wv_from_bin
```

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

Loaded vocab size 3000000

**Note:** If you are receiving out of memory issues on your local machine, try closing other applications to free more memory on your device. You may want to try restarting your machine so that you can free up extra memory. Then immediately run the jupyter notebook and see if you can load the word vectors properly. If you still have problems with loading the embeddings onto your local machine after this, please follow the Piazza instructions, as how to run remotely on Stanford Farmshare machines.

#### 1.3.1 Reducing dimensionality of Word2Vec Word Embeddings

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

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

```
In [15]: def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd', 'ecuador', 'e
         """ Put the word2vec vectors into a matrix M.
         Param:
             wv_from_bin: KeyedVectors object; the 3 million word2vec vectors loaded f
```

*Return:*

*M: numpy matrix shape (num words, 300) containing the vectors  
word2Ind: dictionary mapping each word to its row number in M*

```
"""
import random
words = list(wv_from_bin.vocab.keys())
print("Shuffling words ...")
random.shuffle(words)
words = words[:10000]
print("Putting %i words into word2Ind and matrix M..." % len(words))
word2Ind = {}
M = []
curInd = 0
for w in words:
    try:
        M.append(wv_from_bin.word_vec(w))
        word2Ind[w] = curInd
        curInd += 1
    except KeyError:
        continue
for w in required_words:
    try:
        M.append(wv_from_bin.word_vec(w))
        word2Ind[w] = curInd
        curInd += 1
    except KeyError:
        continue
M = np.stack(M)
print("Done.")
return M, word2Ind
```

```
In [16]: # -----
# 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.
```

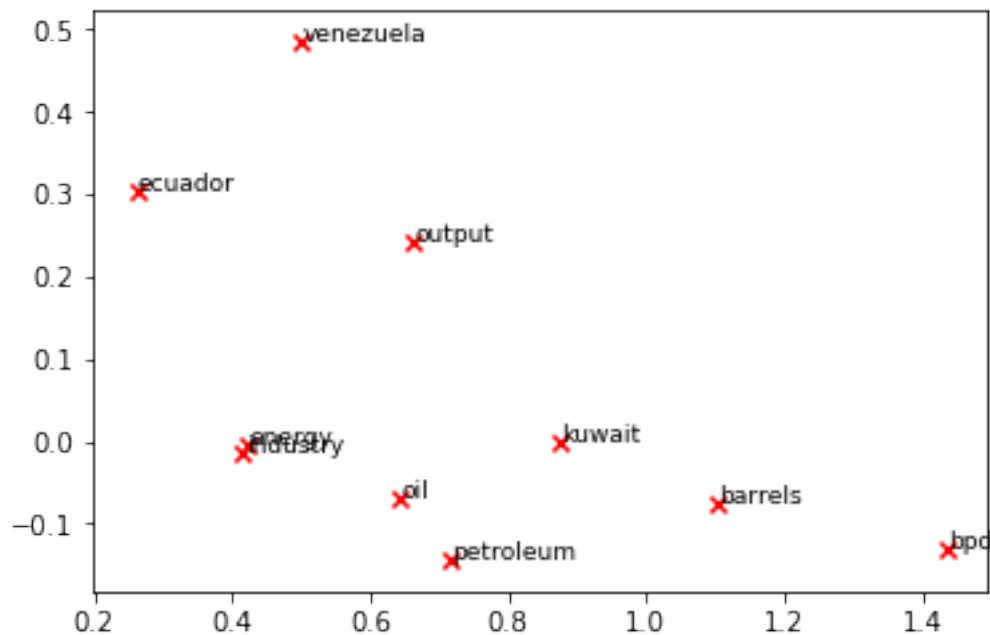
### 1.3.2 Question 2.1: Word2Vec Plot Analysis [written] (4 points)

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



**Write your answer here.** We observe some clusters : - {energy, industry} are clustered together - {venezuela, ecuador} seem nearby in 2-d space

Overall, for the 10 words above, the co-occurrence word embeddings seem more clustered together compared to the word2vec word embeddings.

The co-occurrence word embeddings were created from the reuters corpus for "crude" . The read\_corpus function used to build the co-occurrence embeddings pulls out only articles from the "crude" (i.e. news articles about oil, gas, etc.) category.

For the prediction based word vectors, we are using embeddings produced by word2vec. These vectors have been trained on a corpus of documents from a large number of categories (not just "crude" category). As such, the semantic meaning of a word is much more diverse.

### 1.3.3 Cosine Similarity

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

We can think of n-dimensional vectors as points in n-dimensional space. If we take this perspective L1 and L2 Distances help quantify the amount of space "we must travel" to get between

these two points. Another approach is to examine the angle between two vectors. From trigonometry we know that:

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

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

### 1.3.4 Question 2.2: Polysemous Words (2 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 [18]: # -----
         # Write your polysemous word exploration code here.

         wv_from_bin.most_similar("left")

         # -----
```

```
Out[18]: [('leaving', 0.6707000732421875),
          ('leave', 0.525093138217926),
          ('leaves', 0.5228644609451294),
          ('returned', 0.5059226751327515),
          ('right', 0.49213993549346924),
          ('departed', 0.49109700322151184),
          ('limped', 0.48599502444267273),
          ('went', 0.4719873070716858),
          ('remaining', 0.4650370478630066),
          ('empty', 0.4546155333518982)]
```

**Write your answer here.** "left" has both "leave" and "right" in the top 10. This is because "left" is a semous word. It can mean the act of leaving (hence presence of word "leave") and it also means a direction opposite to "right"

Many of the polysemous words did not have related words from both meanings in the Top-10. I believe this is because one of the meanings is dominant. Here we are looking at the top-10 so words from both meanings did not make it to the top 10 list.

### 1.3.5 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  $1 - \text{Cosine Similarity}$ .

Find three words ( $w_1, w_2, w_3$ ) where  $w_1$  and  $w_2$  are synonyms and  $w_1$  and  $w_3$  are antonyms, but  $\text{Cosine Distance}(w_1, w_3) < \text{Cosine Distance}(w_1, w_2)$ . For example,  $w_1 = \text{"happy"}$  is closer to  $w_3 = \text{"sad"}$  than to  $w_2 = \text{"cheerful"}$ .

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

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

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

        w1 = "sad"
        w2 = "unhappy"
        w3 = "happy"
        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 sad, unhappy have cosine distance: 0.5842775502560309

Antonyms sad, happy have cosine distance: 0.4645385660405297

**Write your answer here.** The synonyms ("sad", "unhappy") have larger cosine distance between them compared to antonyms ("sad", "happy")

The word2vec embedding we are using here is a prediction based word vector. The words ("sad", "happy") are highly predictable in text since they are often used together even though they are antonyms.

### 1.3.6 Solving Analogies with Word Vectors

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

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

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

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

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

        [('queen', 0.7118192315101624),
         ('monarch', 0.6189675331115723),
```

```
(('princess', 0.5902431011199951),
 ('crown_prince', 0.5499460697174072),
 ('prince', 0.5377321243286133),
 ('kings', 0.5236844420433044),
 ('Queen_Consort', 0.5235946178436279),
 ('queens', 0.5181134343147278),
 ('sultan', 0.5098592638969421),
 ('monarchy', 0.5087411999702454])
```

### 1.3.7 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). 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 [21]: # -----
        # Write your analogy exploration code here.

        pprint.pprint(wv_from_bin.most_similar(positive=['india', 'japanese'], negative=['jap

        # -----

[('indian', 0.5673424005508423),
 ('british', 0.5398949384689331),
 ('american', 0.5162887573242188),
 ('chinese', 0.5117965340614319),
 ('swedish', 0.5106680393218994),
 ('canadian', 0.5039363503456116),
 ('pakistani', 0.5012112855911255),
 ('german', 0.5010029077529907),
 ('australia', 0.4979511499404907),
 ('canada', 0.49694547057151794)]
```

**Write your answer here.**    japan:japanese :: india:indian

### 1.3.8 Question 2.5: Incorrect Analogy [code + written] (1 point)

Find an example of analogy that does *not* hold according to these vectors. In your solution, state the intended analogy in the form  $x:y :: a:b$ , and state the (incorrect) value of  $b$  according to the word vectors.

```
In [22]: # -----
        # Write your incorrect analogy exploration code here.

        pprint.pprint(wv_from_bin.most_similar(positive=['indian', 'japan'], negative=['japan
```

```
# -----
[('pakistan', 0.5793417692184448),
 ('india', 0.5646401643753052),
 ('mexico', 0.5408872365951538),
 ('america', 0.5384131073951721),
 ('srilanka', 0.5344810485839844),
 ('sonia', 0.5213700532913208),
 ('indians', 0.5208514928817749),
 ('africa', 0.5049134492874146),
 ('oklahoma', 0.5044816136360168),
 ('rahul', 0.5028336644172668)]
```

**Write your answer here.** The correct analogy should be: japanese:japan :: indian:india  
india is 2nd in the list :)

### 1.3.9 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". What do you find in the top 10?

```
In [23]: # Run this cell
# Here `positive` indicates the list of words to be similar to and `negative` indicates
# 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', 0.552264392375946),
 ('manageress', 0.49151355028152466),
 ('exec', 0.45940810441970825),
 ('Manageress', 0.45598432421684265),
 ('receptionist', 0.4474116563796997),
 ('Jane_Danson', 0.44480547308921814),
 ('Fiz_Jennie_McAlpine', 0.44275766611099243),
 ('Coronation_Street_actress', 0.44275569915771484),
 ('supremo', 0.4409852623939514),
 ('coworker', 0.4398624897003174)]

[('supremo', 0.6097398400306702),
 ('MOTHERWELL_boss', 0.5489562153816223),
 ('CARETAKER_boss', 0.5375303626060486),
 ('Bully_Wee_boss', 0.533397376537323),
```

```
( 'YEOVIL_Town_boss', 0.5321705341339111),
( 'head_honcho', 0.5281979441642761),
( 'manager_Stan_Ternent', 0.5259714722633362),
( 'Viv_Busby', 0.5256163477897644),
( 'striker_Gabby_Agbonlahor', 0.5250812768936157),
( 'BARN斯LEY_boss', 0.5238943696022034)]
```

**Write your answer here.**

- a. "bosses".
- b. "supremo".

"supremo" is #9 in the first list, but #1 in the second list. This reflects bias.

### 1.3.10 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. Please briefly explain the example of bias that you discover.

```
In [24]: # -----
         # Write your bias exploration code here.

         pprint.pprint(wv_from_bin.most_similar(positive=["white", "good"], negative=["black"]))

         # -----

[('nice', 0.6383841037750244),
 ('great', 0.6161423921585083),
 ('terrific', 0.614184558391571),
 ('bad', 0.5943925380706787),
 ('excellent', 0.5763711929321289),
 ('decent', 0.5739346742630005),
 ('fantastic', 0.5718001127243042),
 ('perfect', 0.5275566577911377),
 ('better', 0.5265002250671387),
 ('wonderful', 0.5130241513252258)]
```

**Write your answer here.** The above are the top-10 terms that are most similar to "white" and "good" and most dissimilar to "black".

There seems to be bias exhibited by the vectors. 9 out of the 10 words represent personality traits. Also the word "bad" is #4 in the list which seems counter-intuitive.

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

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

**Write your answer here.** These reflect inherent biases in our writing.

## 2 Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.
2. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
3. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.
4. Once you've rerun everything, select File -> Download as -> PDF via LaTeX
5. Look at the PDF file and make sure all your solutions are there, displayed correctly. The PDF is the only thing your graders will see!
6. Submit your PDF on Gradescope.