

Name: Nayan Vijay Raghatate

Assignment 1: Exploring Word Vectors (25 Points)

Due 11:am, Thurs Sept. 14

Before you start, make sure you read the README.txt in the same directory as this notebook for important setup information. A lot of code is provided in this notebook, and we highly encourage you to read and understand it as part of the learning :)

If you aren't super familiar with Python, Numpy, or Matplotlib, we recommend you check out the CS231N Python/Numpy [tutorial \(https://cs231n.github.io/python-numpy-tutorial/\)](https://cs231n.github.io/python-numpy-tutorial/).

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

In [2]:

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

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

#

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](https://en.wikipedia.org/wiki/Word_embedding) (https://en.wikipedia.org/wiki/Word_embedding) states, "*conceptually it involves a mathematical embedding from a space with one dimension per word to a continuous vector space with a much lower dimension*".

Part 1: Count-Based Word Vectors (10 points)

Most word vector models start from the following idea:

You shall know a word by the company it keeps ([Firth, J. R. 1957:11](https://en.wikipedia.org/wiki/John_Rupert_Firth)
(https://en.wikipedia.org/wiki/John_Rupert_Firth))

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](http://web.stanford.edu/class/cs124/lec/vectorsemantics.video.pdf) (<http://web.stanford.edu/class/cs124/lec/vectorsemantics.video.pdf>) or [here](https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285) (<https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285>)).

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 $n=1$:

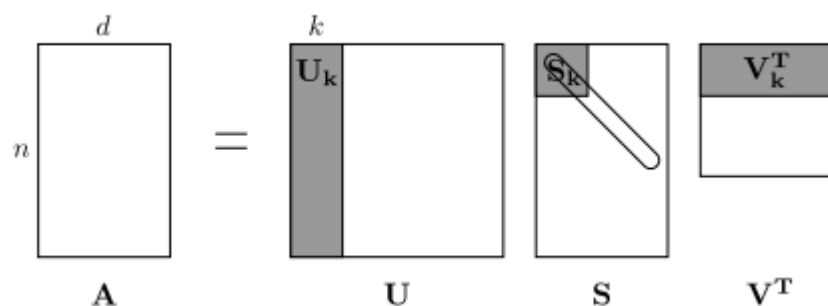
Document 1: "all that glitters is not gold"

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 \(https://davetang.org/file/Singular_Value_Decomposition_Tutorial.pdf\)](https://davetang.org/file/Singular_Value_Decomposition_Tutorial.pdf). If you want to learn more thoroughly about PCA or SVD, feel free to check out lectures [7 \(https://web.stanford.edu/class/cs168/1/17.pdf\)](https://web.stanford.edu/class/cs168/1/17.pdf), [8 \(http://theory.stanford.edu/~tim/s15/1/18.pdf\)](http://theory.stanford.edu/~tim/s15/1/18.pdf), and [9 \(https://web.stanford.edu/class/cs168/1/19.pdf\)](https://web.stanford.edu/class/cs168/1/19.pdf) 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](https://en.wikipedia.org/wiki/Truncated_SVD) (https://en.wikipedia.org/wiki/Singular_value_decomposition#Truncated_SVD)* — then there are reasonably available techniques to compute these 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> (<https://www.nltk.org/book/ch02.html>). We provide a `read_corpus` function below that pulls out only articles from the "crude" (i.e. news articles about oil, gas, etc.) category. The function also adds START and END tokens to each of the documents, and lowercases words. You do **not** have perform any other kind of pre-processing.

```
In [0]: 1 def read_corpus(category="crude"):
2         """ Read files from the specified Reuter's category.
3         Params:
4             category (string): category name
5         Return:
6             list of lists, with words from each of the processed files
7         """
8         files = reuters.fileids(category)
9         return [[START_TOKEN] + [w.lower() for w in list(reuters.words(f))] +
10
```

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

```
In [4]: 1 reuters_corpus = read_corpus()
        2 pprint.pprint(reuters_corpus[:3], compact=True, width=100)
```

[['<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand', 'downwards', 'the',
 'ministry', 'of', 'international', 'trade', 'and', 'industry', '(', 'miti', ')', 'will', 'revise',
 'its', 'long', '-', 'term', 'energy', 'supply', '/', 'demand', 'outlook', 'by', 'august', 'to',
 'meet', 'a', 'forecast', 'downtrend', 'in', 'japanese', 'energy', 'demand', ',', 'ministry',
 'officials', 'said', '.', 'miti', 'is', 'expected', 'to', 'lower', 'the', 'projection', 'for',
 'primary', 'energy', 'supplies', 'in', 'the', 'year', '2000', 'to', '550', 'mln', 'kilolitres',
 '(', 'kl', ')', 'from', '600', 'mln', ',', 'they', 'said', '.', 'the', 'decision', 'follows',
 'the', 'emergence', 'of', 'structural', 'changes', 'in', 'japanese', 'industry', 'following',
 'the', 'rise', 'in', 'the', 'value', 'of', 'the', 'yen', 'and', 'a', 'decline', 'in', 'domestic',
 'electric', 'power', 'demand', '.', 'miti', 'is', 'planning', 'to', 'work', 'out', 'a', 'revised',
 'energy', 'supply', '/', 'demand', 'outlook', 'through', 'deliberations', 'of', 'committee',
 'meetings', 'of', 'the', 'agency', 'of', 'natural', 'resources', 'and', 'energy', ',', 'the',
 'officials', 'said', '.', 'they', 'said', 'miti', 'will', 'also', 'review', 'the', 'breakdown',
 'of', 'energy', 'supply', 'sources', ',', 'including', 'oil', ',', 'nuclear', ',', 'coal', 'and',
 'natural', 'gas', '.', 'nuclear', 'energy', 'provided', 'the', 'bulk', 'of', 'japan', '"', 's',
 'electric', 'power', 'in', 'the', 'fiscal', 'year', 'ended', 'march', '31', ',', 'supplying',
 'an', 'estimated', '27', 'pct', 'on', 'a', 'kilowatt', '/', 'hour', 'basis', ',', 'followed',
 'by', 'oil', '(', '23', 'pct', ')', 'and', 'liquefied', 'natural', 'gas', '(', '21', 'pct', ')',
 'they', 'noted', '.', '<END>'],
 ['<START>', 'energy', '/', 'u', '.', 's', '.', 'petrochemical', 'industry', 'cheap', 'oil',
 'feedstocks', ',', 'the', 'weakened', 'u', '.', 's', '.', 'dollar', 'and', 'a', 'plant',
 'utilization', 'rate', 'approaching', '90', 'pct', 'will', 'propel', 'the', 'streamlined', 'u',
 '.', 's', '.', 'petrochemical', 'industry', 'to', 'record', 'profits', 'this', 'year', ',',
 'with', 'growth', 'expected', 'through', 'at', 'least', '1990', ',', 'major', 'company',
 'executives', 'predicted', '.', 'this', 'bullish', 'outlook', 'for', 'chemical', 'manufacturing',
 'and', 'an', 'industrywide', 'move', 'to', 'shed', 'unrelated', 'businesses', 'has', 'prompted',
 'gaf', 'corp', '&', 'lt', ';', 'gaf', '>', 'privately', '-', 'held', 'cain', 'chemical', 'inc',
 ',', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions', 'of', 'petrochemical',
 'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&', 'lt', ';', 'ash',
]

'>', 'the', 'kentucky', '-', 'based', 'oil', 'refiner', 'and', 'marketer',
, 'are', 'also',
'shopping', 'for', 'money', '-', 'making', 'petrochemical', 'businesses',
'to', 'buy', '.', 'i', 'see', 'us', 'poised', 'at', 'the', 'threshold', 'of', 'a', 'golden',
'period', ', ', 'said',
'paul', 'oreffice', ', ', 'chairman', 'of', 'giant', 'dow', 'chemical', 'c
o', '&', 'lt', ';',
'dow', '>', 'adding', ', ', 'there', 's', 'no', 'major', 'plant',
'capacity', 'being',
'added', 'around', 'the', 'world', 'now', '.', 'the', 'whole', 'game', 'i
s', 'bringing', 'out',
'new', 'products', 'and', 'improving', 'the', 'old', 'ones', '."', 'analyst
s', 'say', 'the',
'chemical', 'industry', 's', 'biggest', 'customers', ', ', 'automobil
e', 'manufacturers',
'and', 'home', 'builders', 'that', 'use', 'a', 'lot', 'of', 'paints', 'an
d', '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', 'hik
e', 'product', 'prices',
'by', '30', 'to', '40', 'pct', 'this', 'year', ', ', 'said', 'john', 'doshe
r', ', ', 'managing',
'director', 'of', 'pace', 'consultants', 'inc', 'of', 'houston', '.', 'dema
nd', 'for', 'some',
'products', 'such', 'as', 'styrene', 'could', 'push', 'profit', 'margins',
'up', 'by', 'as',
'much', 'as', '300', 'pct', ', ', 'he', 'said', '.', 'oreffice', ', ', 'speak
ing', 'at', 'a',
'meeting', 'of', 'chemical', 'engineers', 'in', 'houston', ', ', 'said', 'do
w', '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', 'an
d', 'chemical',
'exports', 'were', 'adversely', 'affected', 'by', 'the', 'strong', 'u',
, 's', 'dollar',
, 'dow', 'had', 'profits', 'of', '58', 'mln', 'dlrs', '.', 'i', 'be
lieve', 'the',
'entire', 'chemical', 'industry', 'is', 'headed', 'for', 'a', 'record', 'ye
ar', 'or', 'close',
'to', 'it', ', ', 'oreffice', 'said', '.', 'gaf', 'chairman', 'samuel', 'he
yman', 'estimated',
'that', 'the', 'u', '.', 's', '.', 'chemical', 'industry', 'would', 'repor
t', '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', 'comb

ination', 'of', 'luck',
'and', 'planning', ',', 'said', 'pace', '"', 's', 'john', 'dosher', '.', 'd
osher', 'said', 'last',
'year', '"', 's', 'fall', 'in', 'oil', 'prices', 'made', 'feedstocks', 'dra
matically', 'cheaper',
'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was', 'we
akening', 'against',
'foreign', 'currencies', '.', 'that', 'helped', 'boost', 'u', '.', 's',
'.', 'chemical',
'exports', '.', 'also', 'helping', 'to', 'bring', 'supply', 'and', 'deman
d', 'into', 'balance',
'has', 'been', 'the', 'gradual', 'market', 'absorption', 'of', 'the', 'extr
a', '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', 'rest
ructuring', 'touched',
'off', 'a', 'flurry', 'of', 'friendly', 'and', 'hostile', 'takeover', 'atte
mpts', '.', 'gaf', ',',
'which', 'made', 'an', 'unsuccessful', 'attempt', 'in', '1985', 'to', 'acqu
ire', 'union',
'carbide', 'corp', '&', 'lt', ';', 'uk', '>', 'recently', 'offered', 'thre
e', 'billion', 'dlrs',
'for', 'borg', 'warner', 'corp', '&', 'lt', ';', 'bor', '>', 'a', 'chicag
o', 'manufacturer',
'of', 'plastics', 'and', 'chemicals', '.', 'another', 'industry', 'powerhou
se', ',', 'w', '.',
'r', '.', 'grace', '&', 'lt', ';', 'gra', '>', 'has', 'divested', 'its', 'r
etailing', ',',
'restaurant', 'and', 'fertilizer', 'businesses', 'to', 'raise', 'cash', 'fo
r', 'chemical',
'acquisitions', '.', 'but', 'some', 'experts', 'worry', 'that', 'the', 'che
mical', 'industry',
'may', 'be', 'headed', 'for', 'trouble', 'if', 'companies', 'continue', 'tu
rning', 'their',
'back', 'on', 'the', 'manufacturing', 'of', 'staple', 'petrochemical', 'com
modities', ',', 'such',
'as', 'ethylene', ',', 'in', 'favor', 'of', 'more', 'profitable', 'specialt
y', 'chemicals',
'that', 'are', 'custom', '-', 'designed', 'for', 'a', 'small', 'group', 'o
f', 'buyers', '.', '"',
'companies', 'like', 'dupont', '&', 'lt', ';', 'dd', '>', 'and', 'monsant
o', 'co', '&', 'lt', ';',
'mtc', '>', 'spent', 'the', 'past', 'two', 'or', 'three', 'years', 'tryin
g', 'to', 'get', 'out',
'of', 'the', 'commodity', 'chemical', 'business', 'in', 'reaction', 'to',
'how', 'badly', 'the',
'market', 'had', 'deteriorated', ',', '"', 'dosher', 'said', '.', '"', 'but',
'i', 'think', 'they',
'will', 'eventually', 'kill', 'the', 'margins', 'on', 'the', 'profitable',
'chemicals', 'in',

'the', 'niche', 'market', '.', 'some', 'top', 'chemical', 'executives', 'share', 'the',
'concern', '.', 'the', 'challenge', 'for', 'our', 'industry', 'is', 'to', 'keep', 'from',
'getting', 'carried', 'away', 'and', 'repeating', 'past', 'mistakes', 'gaf', 's',
'heyman', 'cautioned', '.', 'the', 'shift', 'from', 'commodity', 'chemicals', 'may', 'be',
'ill', '-', 'advised', '.', 'specialty', 'businesses', 'do', 'not', 'stay', 'special', 'long',
'.', 'houston', '-', 'based', 'cain', 'chemical', 'created', 'this', 'month', 'by', 'the',
'sterling', 'investment', 'banking', 'group', 'believes', 'it', 'can', 'generate', '700',
'mln', 'dlrs', 'in', 'annual', 'sales', 'by', 'bucking', 'the', 'industry', 'trend', '.',
'chairman', 'gordon', 'cain', 'who', 'previously', 'led', 'a', 'leveraged', 'buyout', 'of',
'dupont', 's', 'conoco', 'inc', 's', 'chemical', 'business', 'has', 'spent', '1',
'.', '1', 'billion', 'dlrs', 'since', 'january', 'to', 'buy', 'seven', 'petrochemical', 'plants',
'along', 'the', 'texas', 'gulf', 'coast', '.', 'the', 'plants', 'produce', 'only', 'basic',
'commodity', 'petrochemicals', 'that', 'are', 'the', 'building', 'blocks', 'of', 'specialty',
'products', '.', 'this', 'kind', 'of', 'commodity', 'chemical', 'business', 'will', 'never',
'be', 'a', 'glamorous', 'high', '-', 'margin', 'business', 'cain', 'said', 'adding', 'that', 'demand', 'is', 'expected', 'to', 'grow', 'by', 'about', 'three', 'pct',
'annually', '.', 'garo', 'armen', 'an', 'analyst', 'with', 'dean', 'witter', 'reynolds', 'said', 'chemical', 'makers', 'have', 'also', 'benefitted', 'by', 'increasing', 'demand', 'for',
'plastics', 'as', 'prices', 'become', 'more', 'competitive', 'with', 'aluminum', 'wood',
'and', 'steel', 'products', '.', 'armen', 'estimated', 'the', 'upturn', 'in', 'the', 'chemical',
'business', 'could', 'last', 'as', 'long', 'as', 'four', 'or', 'five', 'years', 'provided',
'the', 'u', 's', 'economy', 'continues', 'its', 'modest', 'rate', 'of', 'growth', '.',
'<END>'],
['<START>', 'turkey', 'calls', 'for', 'dialogue', 'to', 'solve', 'dispute', 'turkey', 'said',
'today', 'its', 'disputes', 'with', 'greece', 'including', 'rights', 'on', 'the',
'continental', 'shelf', 'in', 'the', 'aegean', 'sea', 'should', 'be', 'solved', 'through',
'negotiations', '.', 'a', 'foreign', 'ministry', 'statement', 'said', 'the', 'latest', 'crisis',
'between', 'the', 'two', 'nato', 'members', 'stemmed', 'from', 'the', 'continental', 'shelf',
'dispute', 'and', 'an', 'agreement', 'on', 'this', 'issue', 'would', 'effect', 'the', 'security',

```

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

```

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 \(https://coderwall.com/p/rcmaea/flatten-a-list-of-lists-in-one-line-in-python\)](https://coderwall.com/p/rcmaea/flatten-a-list-of-lists-in-one-line-in-python) may be useful to flatten a list of lists. If you're not familiar with Python list comprehensions in general, here's [more information \(https://python-3-patterns-idioms-test.readthedocs.io/en/latest/Comprehensions.html\)](https://python-3-patterns-idioms-test.readthedocs.io/en/latest/Comprehensions.html).

You may find it useful to use [Python sets \(https://www.w3schools.com/python/python_sets.asp\)](https://www.w3schools.com/python/python_sets.asp) to remove duplicate words.

```
In [0]: 1 def distinct_words(corpus):
2         """ Determine a list of distinct words for the corpus.
3         Params:
4             corpus (list of list of strings): corpus of documents
5         Return:
6             corpus_words (list of strings): list of distinct words across
7             num_corpus_words (integer): number of distinct words across th
8         """
9         corpus_words = []
10        num_corpus_words = -1
11
12        # -----
13        # Write your implementation here.
14
15        temp_words = set()
16        for sentence in corpus:
17            for word in sentence:
18                temp_words.add(word)
19        corpus_words = sorted(list(temp_words))
20        num_corpus_words = len(corpus_words)
21
22        # -----
23
24        return corpus_words, num_corpus_words
```

```
In [6]: 1 # -----
2 # Run this sanity check
3 # Note that this not an exhaustive check for correctness.
4 # -----
5
6 # Define toy corpus
7 test_corpus = ["START All that glitters isn't gold END".split(" "), "START
8 test_corpus_words, num_corpus_words = distinct_words(test_corpus)
9
10 # Correct answers
11 ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that", "
12 ans_num_corpus_words = len(ans_test_corpus_words)
13
14 # Test correct number of words
15 assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of dis
16
17 # Test correct words
18 assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_wor
19
20 # Print Success
21 print("-" * 80)
22 print("Passed All Tests!")
23 print("-" * 80)
```

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

Question 1.2: Implement `compute_co_occurrence_matrix` [code] (3 points)

Write a method that constructs a co-occurrence matrix for a certain window-size n (with a default of 4), considering words n before and n after the word in the center of the window. Here, we start to use `numpy` (`np`) to represent vectors, matrices, and tensors. If you're not familiar with NumPy, there's a NumPy tutorial in the second half of this cs231n [Python NumPy tutorial](http://cs231n.github.io/python-numpy-tutorial/) (<http://cs231n.github.io/python-numpy-tutorial/>).

In [0]:

```
1 from collections import defaultdict
2
3 def compute_co_occurrence_matrix(corpus, window_size=4):
4     """ Compute co-occurrence matrix for the given corpus and window_size
5
6         Note: Each word in a document should be at the center of a window.
7             number of co-occurring words.
8
9             For example, if we take the document "START All that glitter
10             "All" will co-occur with "START", "that", "glitters", "is",
11
12     Params:
13         corpus (list of list of strings): corpus of documents
14         window_size (int): size of context window
15     Return:
16         M (numpy matrix of shape (number of corpus words, number of co
17             Co-occurrence matrix of word counts.
18             The ordering of the words in the rows/columns should be th
19             word2Ind (dict): dictionary that maps word to index (i.e. row/
20     """
21     words, num_words = distinct_words(corpus)
22     M = None
23     word2Ind = {}
24
25     # -----
26     # Write your implementation here.
27
28     for ind, word in enumerate(words):
29         word2Ind[word] = ind
30
31     M = np.zeros((num_words, num_words))
32     for sentence in corpus:
33         for i in range(len(sentence)):
34             start = i - window_size if i - window_size >= 0 else 0
35             end = i + window_size if i + window_size < len(sentence) else
36             words_in_window = sentence[start:end]
37             temp_count = defaultdict(int)
38             for word in words_in_window:
39                 temp_count[word2Ind[word]] += 1
40
41             for toCompare in temp_count.keys():
42                 for others, value in temp_count.items():
43                     if toCompare != others:
44                         M[toCompare, others] += value
45
46     # -----
47
48     return M, word2Ind
```

In [8]:

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


Passed All Tests!

Question 1.3: Implement `reduce_to_k_dim` [code] (1 point)

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

Note: All of numpy, scipy, and scikit-learn (`sklearn`) provide *some* implementation of SVD, but only scipy and sklearn provide an implementation of Truncated SVD, and only sklearn provides an efficient randomized algorithm for calculating large-scale Truncated SVD. So please use [sklearn.decomposition.TruncatedSVD](https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html) (<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.TruncatedSVD.html>).

```
In [0]: 1 def reduce_to_k_dim(M, k=2):
2         """ Reduce a co-occurrence count matrix of dimensionality (num_corpus_w
3             to a matrix of dimensionality (num_corpus_words, k) using the foll
4             - http://scikit-learn.org/stable/modules/generated/sklearn.dec
5
6         Params:
7             M (numpy matrix of shape (number of corpus words, number of co
8             k (int): embedding size of each word after dimension reduction
9         Return:
10            M_reduced (numpy matrix of shape (number of corpus words, k)):
11                In terms of the SVD from math class, this actually ret
12        """
13        n_iters = 10      # Use this parameter in your call to `TruncatedSVD`
14        M_reduced = None
15        print("Running Truncated SVD over %i words..." % (M.shape[0]))
16
17        # -----
18        # Write your implementation here.
19
20        svd = TruncatedSVD(n_components=k)
21        M_reduced = svd.fit_transform(M)
22
23        # -----
24
25        print("Done.")
26        return M_reduced
```

In [10]:

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

Running Truncated SVD over 10 words...

Done.

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

Question 1.4: Implement `plot_embeddings` [code] (1 point)

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

For this example, you may find it useful to adapt [this code](#)

(<https://www.pythonmembers.club/2018/05/08/matplotlib-scatter-plot-annotate-set-text-at-label-each-point/>). In the future, a good way to make a plot is to look at [the Matplotlib gallery](#) (<https://matplotlib.org/gallery/index.html>), find a plot that looks somewhat like what you want, and adapt the code they give.

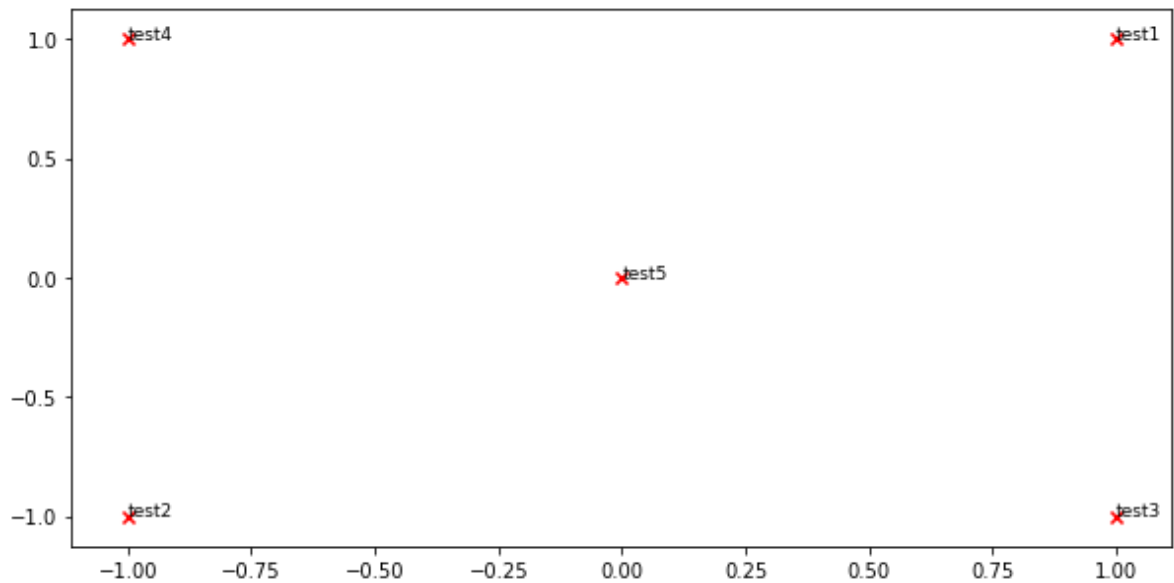
In [0]:

```
1 def plot_embeddings(M_reduced, word2Ind, words):
2     """ Plot in a scatterplot the embeddings of the words specified in the
3         NOTE: do not plot all the words listed in M_reduced / word2Ind.
4         Include a label next to each point.
5
6         Params:
7             M_reduced (numpy matrix of shape (number of unique words in th
8             word2Ind (dict): dictionary that maps word to indices for matr
9             words (list of strings): words whose embeddings we want to vis
10    """
11
12    # -----
13    # Write your implementation here.
14
15    indices = [word2Ind[word] for word in words]
16    x_coors = [M_reduced[index, 0] for index in indices]
17    y_coors = [M_reduced[index, 1] for index in indices]
18
19    for i, word in enumerate(words):
20        x = x_coors[i]
21        y = y_coors[i]
22        plt.scatter(x, y, marker='x', color='red')
23        plt.text(x, y, word, fontsize=9)
24    plt.show()
25
26    # -----
```

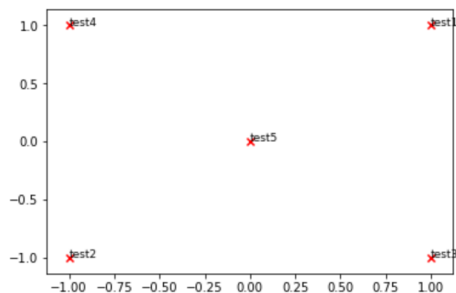
In [12]:

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

Outputted Plot:



Test Plot Solution



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

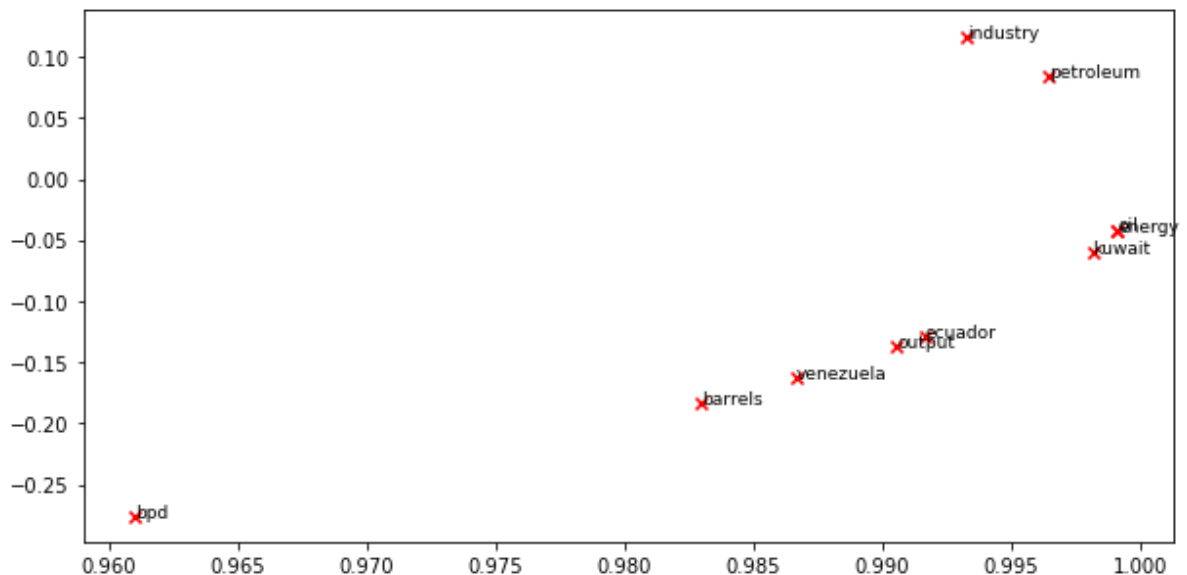
Now we will put together all the parts you have written! We will compute the co-occurrence matrix with fixed window of 4, over the Reuters "crude" corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns $U \cdot S$, so we normalize the returned vectors, so that all the vectors will appear around the unit circle (therefore closeness is directional closeness). **Note:** The line of code below that does the normalizing uses the NumPy concept of *broadcasting*. If you don't know about broadcasting, check out [Computation on Arrays: Broadcasting by Jake VanderPlas](https://jakevdp.github.io/PythonDataScienceHandbook/02.05-computation-on-arrays-broadcasting.html) (<https://jakevdp.github.io/PythonDataScienceHandbook/02.05-computation-on-arrays-broadcasting.html>).

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 [13]:

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

Running Truncated SVD over 8185 words...
Done.



Write your answer here.

What clusters together in 2-dimensional embedding space?

- Group A:
 - industry
 - petroleum
- Group B:
 - energy
 - oil
 - kuwait
- Group C:
 - ecuador
 - output
 - venezuela
 - barrels
- Group D:
 - bpd

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

I think all the place names should be cluster together. Such as Kuwait, Ecuador, Venezuela.

I don't think the performance of this co-occurrence matrix method is doing well.

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](https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf) (<https://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality.pdf>).

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

```
In [0]: 1 def load_word2vec():
2         """ Load Word2Vec Vectors
3         Return:
4             wv_from_bin: All 3 million embeddings, each length 300
5         """
6         import gensim.downloader as api
7         wv_from_bin = api.load("word2vec-google-news-300")
8         vocab = list(wv_from_bin.vocab.keys())
9         print("Loaded vocab size %i" % len(vocab))
10        return wv_from_bin
```

In [5]:

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

/usr/local/lib/python3.6/dist-packages/smart_open/smart_open_lib.py:398: User Warning: This function is deprecated, use smart_open.open instead. See the migration notes for details: https://github.com/RaRe-Technologies/smart_open/blob/master/README.rst#migrating-to-the-new-open-function (https://github.com/RaRe-Technologies/smart_open/blob/master/README.rst#migrating-to-the-new-open-function)

'See the migration notes for details: %s' % _MIGRATION_NOTES_URL

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.

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 [0]: 1 def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd', '
2         """ Put the word2vec vectors into a matrix M.
3         Param:
4             wv_from_bin: KeyedVectors object; the 3 million word2vec vecto
5         Return:
6             M: numpy matrix shape (num words, 300) containing the vectors
7             word2Ind: dictionary mapping each word to its row number in M
8         """
9         import random
10        words = list(wv_from_bin.vocab.keys())
11        print("Shuffling words ...")
12        random.shuffle(words)
13        words = words[:10000]
14        print("Putting %i words into word2Ind and matrix M..." % len(words))
15        word2Ind = {}
16        M = []
17        curInd = 0
18        for w in words:
19            try:
20                M.append(wv_from_bin.word_vec(w))
21                word2Ind[w] = curInd
22                curInd += 1
23            except KeyError:
24                continue
25        for w in required_words:
26            try:
27                M.append(wv_from_bin.word_vec(w))
28                word2Ind[w] = curInd
29                curInd += 1
30            except KeyError:
31                continue
32        M = np.stack(M)
33        print("Done.")
34        return M, word2Ind

```

```

In [9]: 1 # -----
2 # Run Cell to Reduce 300-Dimensional Word Embeddings to k Dimensions
3 # Note: This may take several minutes
4 # -----
5 M, word2Ind = get_matrix_of_vectors(wv_from_bin)
6 M_reduced = reduce_to_k_dim(M, k=2)

```

```

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

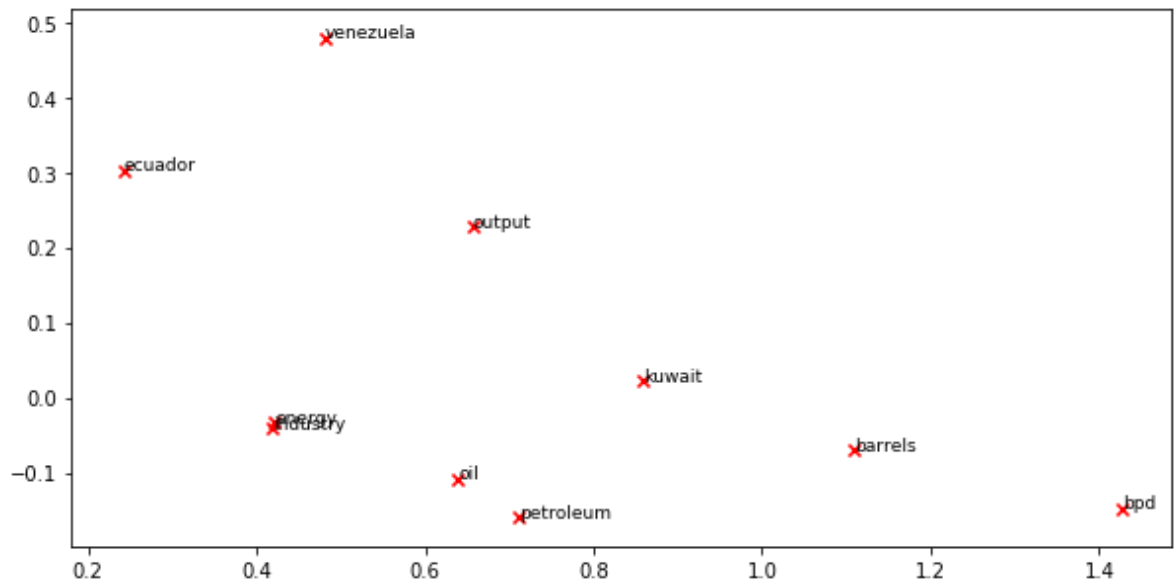
```

Question 2.1: Word2Vec Plot Analysis [written] (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 [12]: 1 words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil',  
2 plot_embeddings(M_reduced, word2Ind, words)
```



Write your answer here.

What clusters together in 2-dimensional embedding space?

Maybe just

- industry
- energy

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

Almost every things that are seperated from each other for about same distance.

Not as I expect that the place name will cluster together.

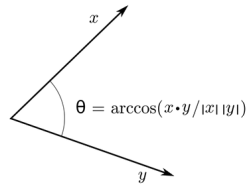
How is the plot different from the one generated earlier from the co-occurrence matrix?

The clustering phenomemon is much less obvious here.

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](https://en.wikipedia.org/wiki/Cosine_similarity) (https://en.wikipedia.org/wiki/Cosine_similarity) s between two vectors p and q is defined as:

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

Question 2.2: Polysemous Words (2 points) [code + written]

Find a [polysemous](https://en.wikipedia.org/wiki/Polysemy) (<https://en.wikipedia.org/wiki/Polysemy>) 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](https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors) (<https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors>).



In [13]:

```
1 # -----
2 # Write your polysemous word exploration code here.
3
4 pprint.pprint(wv_from_bin.most_similar('leaves'))
5 pprint.pprint(wv_from_bin.most_similar('scoop'))
6 pprint.pprint(wv_from_bin.most_similar('bank'))
7
8 # -----
```

/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of issubdtype from `int` to `np.signedinteger` is deprecated. In future, it will be treated as `np.int64 == np.dtype(int).type`.

```
if np.issubdtype(vec.dtype, np.int):
```

```
[('leaving', 0.5886485576629639),
 ('arrives', 0.5362045764923096),
 ('left', 0.5228645205497742),
 ('Leaves', 0.5153512954711914),
 ('leave', 0.5131404399871826),
 ('departs', 0.5107832551002502),
 ('lingers', 0.4776709973812103),
 ('vanishes', 0.4696614742279053),
 ('stalks', 0.46650534868240356),
 ('sends', 0.4623021185398102)]
[('scoops', 0.6741111278533936),
 ('scooped', 0.5712401270866394),
 ('Royal_Wedding_Meltdowns', 0.5508615374565125),
 ('Restrictions_Leash', 0.5447814464569092),
 ('Use_melon_baller', 0.5258612632751465),
 ('scooping', 0.5147117376327515),
 ('news@morehorror.com', 0.5114910006523132),
 ('handed_waffle_cone', 0.49177810549736023),
 ('lowdown', 0.4857146143913269),
 ('techie_breakie', 0.47323939204216003)]
[('banks', 0.7440759539604187),
 ('banking', 0.690161406993866),
 ('Bank', 0.6698698997497559),
 ('lender', 0.6342284679412842),
 ('banker', 0.6092953681945801),
 ('depositors', 0.6031531691551208),
 ('mortgage_lender', 0.5797975659370422),
 ('depositor', 0.5716428160667419),
 ('BofA', 0.5714625120162964),
 ('Citibank', 0.5589520335197449)]
```

Write your answer here.

Why do you think many of the polysemous words you tried didn't work?

I feel that finding a polysemous word where the top-10 most similar words contain related words from both meanings can be challenging due to the limitations of word vectors and also few reasons such as:

1. While representing the words in high dimensional vector space, words with multiple meanings might have their different senses mixed together.
2. Word vectors are often trained on large text corpora, thus less common meanings or senses of polysemous words might not have strong representations due to the dominance of frequently appearing words.
3. A word's meaning can change according to the context and word vectors don't capture context.
4. Size of the training corpus can impact the ability to capture subtle differences between words.

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](https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors) (<https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors>) for further assistance.

In [14]:

```

1  # -----
2  # Write your synonym & antonym exploration code here.
3
4  w1 = "fun"
5  w2 = "interesting"
6  w3 = "boring"
7  w1_w2_dist = wv_from_bin.distance(w1, w2)
8  w1_w3_dist = wv_from_bin.distance(w1, w3)
9
10 print("Synonyms {}, {} have cosine distance: {}".format(w1, w2, w1_w2_dist))
11 print("Antonyms {}, {} have cosine distance: {}".format(w1, w3, w1_w3_dist))
12
13 # -----

```

Synonyms fun, interesting have cosine distance: 0.5479990839958191

Antonyms fun, boring have cosine distance: 0.5349873304367065

/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of issubdtype from `int` to `np.signedinteger` is deprecated. In future, it will be treated as `np.int64 == np.dtype(int).type`.

```
if np.issubdtype(vec.dtype, np.int):
```

Write your answer here.

Why this counter-intuitive (antonyms' distance < synonyms' distance) result may have happened?

1. Words like "fun" and "interesting" may have multiple senses or meanings. In some senses or contexts, they could be closer in meaning, while in others, they may not be. Here "fun" and "interesting" are indicating a quality or event but in other context it may not be same.
2. Word vectors are a simplification of language, and they may not capture all semantic relationships accurately. They might excel at capturing broad semantic similarities but might struggle with subtle nuances like the difference between synonyms and antonyms in specific contexts.
3. Example: "fun" and "interesting" are considered synonyms in many contexts, and their vector representations might reflect this similarity. "fun" and "boring" are antonyms, but the vectors may not capture this antonymy as effectively due to the potential variability in how these words are used in different contexts

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](https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors) (<https://radimrehurek.com/gensim/models/keyedvectors.html#gensim.models.keyedvectors>)

```
In [15]: 1 # Run this cell to answer the analogy -- man : king :: woman : x
          2 pprint.pprint(wv_from_bin.most_similar(positive=['woman', 'king'], negativ
```

```
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning:
Conversion of the second argument of issubdtype from `int` to `np.signedinteg
er` is deprecated. In future, it will be treated as `np.int64 == np.dtype(in
t).type`.
```

```
    if np.issubdtype(vec.dtype, np.int):
```

```
[('queen', 0.7118192911148071),
 ('monarch', 0.6189674139022827),
 ('princess', 0.5902431011199951),
 ('crown_prince', 0.5499460697174072),
 ('prince', 0.5377321243286133),
 ('kings', 0.5236844420433044),
 ('Queen_Consort', 0.5235945582389832),
 ('queens', 0.518113374710083),
 ('sultan', 0.5098593235015869),
 ('monarchy', 0.5087411999702454)]
```

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

Find an example of analogy that holds according to these vectors (i.e. the intended word is ranked top). 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 [16]:

```
1 # -----
2 # Write your analogy exploration code here.
3
4 pprint.pprint(wv_from_bin.most_similar(positive=['fantastic', 'good'], neg
5 pprint.pprint(wv_from_bin.most_similar(positive=['chinese', 'america'], ne
6
7 # -----
```

```
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning:
Conversion of the second argument of issubdtype from `int` to `np.signedinteg
er` is deprecated. In future, it will be treated as `np.int64 == np.dtype(in
t).type`.
```

```
if np.issubdtype(vec.dtype, np.int):
```

```
[('terrific', 0.781560480594635),
 ('great', 0.7748532891273499),
 ('wonderful', 0.7381185293197632),
 ('excellent', 0.70376056432724),
 ('marvelous', 0.6975229382514954),
 ('phenomenal', 0.6722577810287476),
 ('fabulous', 0.6698534488677979),
 ('amazing', 0.6655664443969727),
 ('awesome', 0.651089072227478),
 ('nice', 0.6450716853141785)]
[('american', 0.6357719898223877),
 ('americans', 0.5996809601783752),
 ('australian', 0.5581253170967102),
 ('usa', 0.54793781042099),
 ('indian', 0.5451273918151855),
 ('india', 0.5419946908950806),
 ('europe', 0.5409786701202393),
 ('texas', 0.5398456454277039),
 ('mexico', 0.5391411185264587),
 ('british', 0.5385651588439941)]
```

Write your answer here.

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?

The analogy "man : king :: woman : x" can be stated as:

"man" is to "king" as "woman" is to "queen."

1. In this analogy, we are looking for a word that represents the female counterpart of a royal title, just as "king" represents the male counterpart of a royal title. The word vectors correctly identify "queen" as the most similar word when you provide "woman" and "king" as positive examples and "man" as a negative example.
2. The analogy "fantastic : good :: bad : x" can be stated as: "fantastic" is to "good" as "bad" is to "horrible.".....This analogy holds because it is capturing the relationship between words with similar meanings but different degrees of intensity. "Fantastic" is a stronger, more

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 [17]:

```
1 # -----
2 # Write your incorrect analogy exploration code here.
3
4 pprint.pprint(wv_from_bin.most_similar(positive=['soil', 'food'], negative
5
6 # -----
```

/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning: Conversion of the second argument of issubdtype from `int` to `np.signedinteger` is deprecated. In future, it will be treated as `np.int64 == np.dtype(int).type`.

```
if np.issubdtype(vec.dtype, np.int):
```

```
[('foods', 0.48887646198272705),
 ('cooked_meals', 0.4297821819782257),
 ('meals', 0.40927883982658386),
 ('soils', 0.40443646907806396),
 ('Food', 0.3950459957122803),
 ('staple_foods', 0.39481669664382935),
 ('food_stuffs', 0.39097312092781067),
 ('grub', 0.3828972578048706),
 ('nutritious', 0.380367636680603),
 ('foodstuffs', 0.37955567240715027)]
```

Write your answer here.

State the (incorrect) value of b according to the word vectors

1. The analogy is in the form $x:y::a:b$ i.e if "soil" is to "food" as "plant" is to b
2. In the above case I was expecting that "plant" would be related to some term related to the growth or production of food. But, However, the word vectors incorrectly suggest "foods" as the most similar word, which is not an appropriate term in this context.

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 [18]:

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

```
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning:
Conversion of the second argument of issubdtype from `int` to `np.signedinteg
er` is deprecated. In future, it will be treated as `np.int64 == np.dtype(in
t).type`.
```

```
if np.issubdtype(vec.dtype, np.int):
```

```
[('bosses', 0.5522644519805908),
 ('manageress', 0.49151360988616943),
 ('exec', 0.459408164024353),
 ('Manageress', 0.45598435401916504),
 ('receptionist', 0.4474116861820221),
 ('Jane_Danson', 0.44480547308921814),
 ('Fiz_Jennie_McAlpine', 0.44275766611099243),
 ('Coronation_Street_actress', 0.44275569915771484),
 ('supremo', 0.4409852921962738),
 ('coworker', 0.4398624897003174)]
```

```
[('supremo', 0.6097397804260254),
 ('MOTHERWELL_boss', 0.5489562153816223),
 ('CARETAKER_boss', 0.5375303626060486),
 ('Bully_Wee_boss', 0.5333974361419678),
 ('YEOVIL_Town_boss', 0.5321705341339111),
 ('head_honcho', 0.5281980037689209),
 ('manager_Stan_Ternent', 0.525971531867981),
 ('Viv_Busby', 0.5256163477897644),
 ('striker_Gabby_Agbonlahor', 0.5250812768936157),
 ('BARNSELEY_boss', 0.5238943099975586)]
```

Write your answer here.

What do you find in the top 10?

1. When we try to find similar words related to "woman" and "boss" keeping "man" as a negative term we find: 'bosses', 'manageress', 'exec', 'receptionist', "Coronation_Street_actress", etc

2. While when we try to find similar words related to "man" and "boss" keeping "women" as a negative term we find: "supremo", "MOTHERWELL_boss", "BARNSELEY_boss", etc
3. This is happening due to gender bias/ gender stereotypes in certain occupations.

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 [19]:

```
1 # -----
2 # Write your bias exploration code here.
3
4 pprint.pprint(wv_from_bin.most_similar(positive=['man', 'waiter'], negative=['woman', 'boss'],
5 print()
6 pprint.pprint(wv_from_bin.most_similar(positive=['king', 'president'], negative=['queen', 'boss'],
7
8 # -----
```

```
/usr/local/lib/python3.6/dist-packages/gensim/matutils.py:737: FutureWarning:
Conversion of the second argument of issubdtype from `int` to `np.signedinteger`
is deprecated. In future, it will be treated as `np.int64 == np.dtype(int).type`.
```

```
if np.issubdtype(vec.dtype, np.int):
```

```
[('maitre_d', 0.5932615995407104),
 ('barman', 0.5798264741897583),
 ('counterman', 0.5765155553817749),
 ('maitre_d', 0.5741099119186401),
 ('waiters', 0.5707054138183594),
 ('busboy', 0.5484530925750732),
 ('bartender', 0.5442490577697754),
 ('headwaiter', 0.5392516851425171),
 ('maitre_d', 0.534460186958313),
 ('doorman', 0.5177234411239624)]

[('President', 0.6191233396530151),
 ('CEO', 0.5661985874176025),
 ('chairman', 0.5635077953338623),
 ('chief_executive', 0.5602068901062012),
 ('vice_president', 0.5367822647094727),
 ('president', 0.49823909997940063),
 ('founder', 0.4955834746360779),
 ('Chairman', 0.4952137768268585),
 ('COO', 0.4799763858318329),
 ('Vice_President', 0.4791256785392761)]
```

Write your answer here.

Briefly explain the example of bias that you discover

1. In the first case the expected output was "waitress" but it didn't give the correct answer due to the gender bias same as above problem because these words are related to roles in the restaurant and service industry which eventually reflect a gender bias by suggesting that

women are associated with roles like 'waiter' in a similar way that men are associated with roles like 'maitre_d' or 'barman.'

2. In second case also it reflects gender bias the analogy "king : president :: queen : x" reveals a bias as well. The top similar words to "queen" and "president" (with "king" as the negative term) include terms like 'CEO,' 'chairman,' and 'chief_executive.' This bias suggests that leadership roles in business and politics are more closely associated with male terms (e.g., 'CEO,' 'chairman') while female terms (e.g., 'queen') are not as strongly related to such roles.

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.

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

1. Word vectors are often trained on large text corpora, thus less common meanings or senses of polysemous words might not have strong representations due to the dominance of frequently appearing words and also words with multiple meanings might have their different senses mixed together, thus it can have word ambiguity.
2. The gender bias due to the presence of biases in such as news articles or historical documents due to the time period or cultural norms they represent which are called as stereotype.

Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.
2. Please make sure to have entered your SUNET ID above.
3. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
4. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.
5. Once you've rerun everything, select File -> Download as -> PDF via LaTeX
6. 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!
7. Submit your PDF on Gradescope.