# exploring_word_vectors

April 8, 2020

## 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.

```
[33]: # 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)
      # ----------------
```

```
[nltk_data] Downloading package reuters to
```

```
[nltk_data]      C:\Users\z8010\AppData\Roaming\nltk_data…
[nltk_data]    Package reuters is already up-to-date!
```

## 1.1   Please Write Your SUNet ID Here: folk19

# 2   Word Vectors

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

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

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

## 2.1   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).

### 2.1.1   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 $w_{i-n} \ldots w_{i-1}$ and $w_{i+1} \ldots w_{i+n}$. 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.

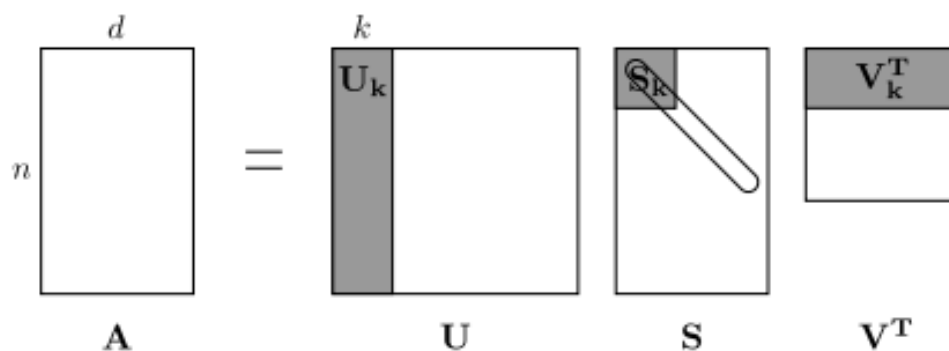**Example: Co-Occurrence with Fixed Window of n=1:**

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 thise 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.

### 2.1.2 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 perform any other kind of pre-processing.

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

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

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

```
[['<START>', 'japan', 'to', 'revise', 'long', '-', 'term', 'energy', 'demand',
'downwards', 'the',
  'ministry', 'of', 'international', 'trade', 'and', 'industry', '(', 'miti',
')', 'will', 'revise',
  'its', 'long', '-', 'term', 'energy', 'supply', '/', 'demand', 'outlook',
'by', 'august', 'to',
  'meet', 'a', 'forecast', 'downtrend', 'in', 'japanese', 'energy', 'demand',
',', 'ministry',
  'officials', 'said', '.', 'miti', 'is', 'expected', 'to', 'lower', 'the',
'projection', 'for',
  'primary', 'energy', 'supplies', 'in', 'the', 'year', '2000', 'to', '550',
'mln', 'kilolitres',
  '(', 'kl', ')', 'from', '600', 'mln', ',', 'they', 'said', '.', 'the',
'decision', 'follows',
```

4

'the', 'emergence', 'of', 'structural', 'changes', 'in', 'japanese',
'industry', 'following',
'the', 'rise', 'in', 'the', 'value', 'of', 'the', 'yen', 'and', 'a',
'decline', 'in', 'domestic',
'electric', 'power', 'demand', '.', 'miti', 'is', 'planning', 'to', 'work',
'out', 'a', 'revised',
'energy', 'supply', '/', 'demand', 'outlook', 'through', 'deliberations',
'of', 'committee',
'meetings', 'of', 'the', 'agency', 'of', 'natural', 'resources', 'and',
'energy', ',', 'the',
'officials', 'said', '.', 'they', 'said', 'miti', 'will', 'also', 'review',
'the', 'breakdown',
'of', 'energy', 'supply', 'sources', ',', 'including', 'oil', ',', 'nuclear',
',', 'coal', 'and',
'natural', 'gas', '.', 'nuclear', 'energy', 'provided', 'the', 'bulk', 'of',
'japan', "'", 's',
'electric', 'power', 'in', 'the', 'fiscal', 'year', 'ended', 'march', '31',
',', 'supplying',
'an', 'estimated', '27', 'pct', 'on', 'a', 'kilowatt', '/', 'hour', 'basis',
',', 'followed',
'by', 'oil', '(', '23', 'pct', ')', 'and', 'liquefied', 'natural', 'gas', '(',
'21', 'pct', '),',
'they', 'noted', '.', '<END>'],
['<START>', 'energy', '/', 'u', '.', 's', '.', 'petrochemical', 'industry',
'cheap', 'oil',
'feedstocks', ',', 'the', 'weakened', 'u', '.', 's', '.', 'dollar', 'and',
'a', 'plant',
'utilization', 'rate', 'approaching', '90', 'pct', 'will', 'propel', 'the',
'streamlined', 'u',
'.', 's', '.', 'petrochemical', 'industry', 'to', 'record', 'profits', 'this',
'year', ',',
'with', 'growth', 'expected', 'through', 'at', 'least', '1990', ',', 'major',
'company',
'executives', 'predicted', '.', 'this', 'bullish', 'outlook', 'for',
'chemical', 'manufacturing',
'and', 'an', 'industrywide', 'move', 'to', 'shed', 'unrelated', 'businesses',
'has', 'prompted',
'gaf', 'corp', '&', 'lt', ';', 'gaf', '>,', 'privately', '-', 'held', 'cain',
'chemical', 'inc',
',', 'and', 'other', 'firms', 'to', 'aggressively', 'seek', 'acquisitions',
'of', 'petrochemical',
'plants', '.', 'oil', 'companies', 'such', 'as', 'ashland', 'oil', 'inc', '&',
'lt', ';', 'ash',
'>,', 'the', 'kentucky', '-', 'based', 'oil', 'refiner', 'and', 'marketer',
',', 'are', 'also',
'shopping', 'for', 'money', '-', 'making', 'petrochemical', 'businesses',
'to', 'buy', '.', '"',
'i', 'see', 'us', 'poised', 'at', 'the', 'threshold', 'of', 'a', 'golden',

'period', ',"', 'said',
  'paul', 'oreffice', ',', 'chairman', 'of', 'giant', 'dow', 'chemical', 'co',
'&', 'lt', ';',
  'dow', '>,', 'adding', ',', '"', 'there', "'", 's', 'no', 'major', 'plant',
'capacity', 'being',
  'added', 'around', 'the', 'world', 'now', '.', 'the', 'whole', 'game', 'is',
'bringing', 'out',
  'new', 'products', 'and', 'improving', 'the', 'old', 'ones', '."', 'analysts',
'say', 'the',
  'chemical', 'industry', "'", 's', 'biggest', 'customers', ',', 'automobile',
'manufacturers',
  'and', 'home', 'builders', 'that', 'use', 'a', 'lot', 'of', 'paints', 'and',
'plastics', ',',
  'are', 'expected', 'to', 'buy', 'quantities', 'this', 'year', '.', 'u', '.',
's', '.',
  'petrochemical', 'plants', 'are', 'currently', 'operating', 'at', 'about',
'90', 'pct',
  'capacity', ',', 'reflecting', 'tighter', 'supply', 'that', 'could', 'hike',
'product', 'prices',
  'by', '30', 'to', '40', 'pct', 'this', 'year', ',', 'said', 'john', 'dosher',
',', 'managing',
  'director', 'of', 'pace', 'consultants', 'inc', 'of', 'houston', '.',
'demand', 'for', 'some',
  'products', 'such', 'as', 'styrene', 'could', 'push', 'profit', 'margins',
'up', 'by', 'as',
  'much', 'as', '300', 'pct', ',', 'he', 'said', '.', 'oreffice', ',',
'speaking', 'at', 'a',
  'meeting', 'of', 'chemical', 'engineers', 'in', 'houston', ',', 'said', 'dow',
'would', 'easily',
  'top', 'the', '741', 'mln', 'dlrs', 'it', 'earned', 'last', 'year', 'and',
'predicted', 'it',
  'would', 'have', 'the', 'best', 'year', 'in', 'its', 'history', '.', 'in',
'1985', ',', 'when',
  'oil', 'prices', 'were', 'still', 'above', '25', 'dlrs', 'a', 'barrel', 'and',
'chemical',
  'exports', 'were', 'adversely', 'affected', 'by', 'the', 'strong', 'u', '.',
's', '.', 'dollar',
  ',', 'dow', 'had', 'profits', 'of', '58', 'mln', 'dlrs', '.', '"', 'i',
'believe', 'the',
  'entire', 'chemical', 'industry', 'is', 'headed', 'for', 'a', 'record',
'year', 'or', 'close',
  'to', 'it', ',"', 'oreffice', 'said', '.', 'gaf', 'chairman', 'samuel',
'heyman', 'estimated',
  'that', 'the', 'u', '.', 's', '.', 'chemical', 'industry', 'would', 'report',
'a', '20', 'pct',
  'gain', 'in', 'profits', 'during', '1987', '.', 'last', 'year', ',', 'the',
'domestic',
  'industry', 'earned', 'a', 'total', 'of', '13', 'billion', 'dlrs', ',', 'a',

'54', 'pct', 'leap',
  'from', '1985', '.', 'the', 'turn', 'in', 'the', 'fortunes', 'of', 'the',
'once', '-', 'sickly',
  'chemical', 'industry', 'has', 'been', 'brought', 'about', 'by', 'a',
'combination', 'of', 'luck',
  'and', 'planning', ',', 'said', 'pace', "'", 's', 'john', 'dosher', '.',
'dosher', 'said', 'last',
  'year', "'", 's', 'fall', 'in', 'oil', 'prices', 'made', 'feedstocks',
'dramatically', 'cheaper',
  'and', 'at', 'the', 'same', 'time', 'the', 'american', 'dollar', 'was',
'weakening', 'against',
  'foreign', 'currencies', '.', 'that', 'helped', 'boost', 'u', '.', 's', '.',
'chemical',
  'exports', '.', 'also', 'helping', 'to', 'bring', 'supply', 'and', 'demand',
'into', 'balance',
  'has', 'been', 'the', 'gradual', 'market', 'absorption', 'of', 'the', 'extra',
'chemical',
  'manufacturing', 'capacity', 'created', 'by', 'middle', 'eastern', 'oil',
'producers', 'in',
  'the', 'early', '1980s', '.', 'finally', ',', 'virtually', 'all', 'major',
'u', '.', 's', '.',
  'chemical', 'manufacturers', 'have', 'embarked', 'on', 'an', 'extensive',
'corporate',
  'restructuring', 'program', 'to', 'mothball', 'inefficient', 'plants', ',',
'trim', 'the',
  'payroll', 'and', 'eliminate', 'unrelated', 'businesses', '.', 'the',
'restructuring', 'touched',
  'off', 'a', 'flurry', 'of', 'friendly', 'and', 'hostile', 'takeover',
'attempts', '.', 'gaf', ',',
  'which', 'made', 'an', 'unsuccessful', 'attempt', 'in', '1985', 'to',
'acquire', 'union',
  'carbide', 'corp', '&', 'lt', ';', 'uk', '>,', 'recently', 'offered', 'three',
'billion', 'dlrs',
  'for', 'borg', 'warner', 'corp', '&', 'lt', ';', 'bor', '>,', 'a', 'chicago',
'manufacturer',
  'of', 'plastics', 'and', 'chemicals', '.', 'another', 'industry',
'powerhouse', ',', 'w', '.',
  'r', '.', 'grace', '&', 'lt', ';', 'gra', '>', 'has', 'divested', 'its',
'retailing', ',',
  'restaurant', 'and', 'fertilizer', 'businesses', 'to', 'raise', 'cash', 'for',
'chemical',
  'acquisitions', '.', 'but', 'some', 'experts', 'worry', 'that', 'the',
'chemical', 'industry',
  'may', 'be', 'headed', 'for', 'trouble', 'if', 'companies', 'continue',
'turning', 'their',
  'back', 'on', 'the', 'manufacturing', 'of', 'staple', 'petrochemical',
'commodities', ',', 'such',
  'as', 'ethylene', ',', 'in', 'favor', 'of', 'more', 'profitable', 'specialty',

'chemicals',
  'that', 'are', 'custom', '-', 'designed', 'for', 'a', 'small', 'group', 'of',
'buyers', '.', '"',
  'companies', 'like', 'dupont', '&', 'lt', ';', 'dd', '>', 'and', 'monsanto',
'co', '&', 'lt', ';',
  'mtc', '>', 'spent', 'the', 'past', 'two', 'or', 'three', 'years', 'trying',
'to', 'get', 'out',
  'of', 'the', 'commodity', 'chemical', 'business', 'in', 'reaction', 'to',
'how', 'badly', 'the',
  'market', 'had', 'deteriorated', ',"', 'dosher', 'said', '.', '"', 'but', 'i',
'think', 'they',
  'will', 'eventually', 'kill', 'the', 'margins', 'on', 'the', 'profitable',
'chemicals', 'in',
  'the', 'niche', 'market', '."', 'some', 'top', 'chemical', 'executives',
'share', 'the',
  'concern', '.', '"', 'the', 'challenge', 'for', 'our', 'industry', 'is', 'to',
'keep', 'from',
  'getting', 'carried', 'away', 'and', 'repeating', 'past', 'mistakes', ',"',
'gaf', "'", 's',
  'heyman', 'cautioned', '.', '"', 'the', 'shift', 'from', 'commodity',
'chemicals', 'may', 'be',
  'ill', '-', 'advised', '.', 'specialty', 'businesses', 'do', 'not', 'stay',
'special', 'long',
  '."', 'houston', '-', 'based', 'cain', 'chemical', ',', 'created', 'this',
'month', 'by', 'the',
  'sterling', 'investment', 'banking', 'group', ',', 'believes', 'it', 'can',
'generate', '700',
  'mln', 'dlrs', 'in', 'annual', 'sales', 'by', 'bucking', 'the', 'industry',
'trend', '.',
  'chairman', 'gordon', 'cain', ',', 'who', 'previously', 'led', 'a',
'leveraged', 'buyout', 'of',
  'dupont', "'", 's', 'conoco', 'inc', "'", 's', 'chemical', 'business', ',',
'has', 'spent', '1',
  '.', '1', 'billion', 'dlrs', 'since', 'january', 'to', 'buy', 'seven',
'petrochemical', 'plants',
  'along', 'the', 'texas', 'gulf', 'coast', '.', 'the', 'plants', 'produce',
'only', 'basic',
  'commodity', 'petrochemicals', 'that', 'are', 'the', 'building', 'blocks',
'of', 'specialty',
  'products', '.', '"', 'this', 'kind', 'of', 'commodity', 'chemical',
'business', 'will', 'never',
  'be', 'a', 'glamorous', ',', 'high', '-', 'margin', 'business', ',"', 'cain',
'said', ',',
  'adding', 'that', 'demand', 'is', 'expected', 'to', 'grow', 'by', 'about',
'three', 'pct',
  'annually', '.', 'garo', 'armen', ',', 'an', 'analyst', 'with', 'dean',
'witter', 'reynolds', ',',
  'said', 'chemical', 'makers', 'have', 'also', 'benefitted', 'by',

'increasing', 'demand', 'for',
  'plastics', 'as', 'prices', 'become', 'more', 'competitive', 'with',
'aluminum', ',', 'wood',
  'and', 'steel', 'products', '.', 'armen', 'estimated', 'the', 'upturn', 'in',
'the', 'chemical',
  'business', 'could', 'last', 'as', 'long', 'as', 'four', 'or', 'five',
'years', ',', 'provided',
  'the', 'u', '.', 's', '.', 'economy', 'continues', 'its', 'modest', 'rate',
'of', 'growth', '.',
  '<END>'],
 ['<START>', 'turkey', 'calls', 'for', 'dialogue', 'to', 'solve', 'dispute',
'turkey', 'said',
  'today', 'its', 'disputes', 'with', 'greece', ',', 'including', 'rights',
'on', 'the',
  'continental', 'shelf', 'in', 'the', 'aegean', 'sea', ',', 'should', 'be',
'solved', 'through',
  'negotiations', '.', 'a', 'foreign', 'ministry', 'statement', 'said', 'the',
'latest', 'crisis',
  'between', 'the', 'two', 'nato', 'members', 'stemmed', 'from', 'the',
'continental', 'shelf',
  'dispute', 'and', 'an', 'agreement', 'on', 'this', 'issue', 'would', 'effect',
'the', 'security',
  ',', 'economy', 'and', 'other', 'rights', 'of', 'both', 'countries', '.', '"',
'as', 'the',
  'issue', 'is', 'basicly', 'political', ',', 'a', 'solution', 'can', 'only',
'be', 'found', 'by',
  'bilateral', 'negotiations', ',"', 'the', 'statement', 'said', '.', 'greece',
'has', 'repeatedly',
  'said', 'the', 'issue', 'was', 'legal', 'and', 'could', 'be', 'solved', 'at',
'the',
  'international', 'court', 'of', 'justice', '.', 'the', 'two', 'countries',
'approached', 'armed',
  'confrontation', 'last', 'month', 'after', 'greece', 'announced', 'it',
'planned', 'oil',
  'exploration', 'work', 'in', 'the', 'aegean', 'and', 'turkey', 'said', 'it',
'would', 'also',
  'search', 'for', 'oil', '.', 'a', 'face', '-', 'off', 'was', 'averted',
'when', 'turkey',
  'confined', 'its', 'research', 'to', 'territorrial', 'waters', '.', '"',
'the', 'latest',
  'crises', 'created', 'an', 'historic', 'opportunity', 'to', 'solve', 'the',
'disputes', 'between',
  'the', 'two', 'countries', ',"', 'the', 'foreign', 'ministry', 'statement',
'said', '.', 'turkey',
  '"', 's', 'ambassador', 'in', 'athens', ',', 'nazmi', 'akiman', ',', 'was',
'due', 'to', 'meet',
  'prime', 'minister', 'andreas', 'papandreou', 'today', 'for', 'the', 'greek',
'reply', 'to', 'a',

```
    'message', 'sent', 'last', 'week', 'by', 'turkish', 'prime', 'minister',
'turgut', 'ozal', '.',
    'the', 'contents', 'of', 'the', 'message', 'were', 'not', 'disclosed', '.',
'<END>']]
```

### 2.1.3  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.

```
[36]: def distinct_words(corpus):
          """ Determine a list of distinct words for the corpus.
              Params:
                  corpus (list of list of strings): corpus of documents
              Return:
                  corpus_words (list of strings): list of distinct words across the␣
      ↪corpus, sorted (using python 'sorted' function)
                  num_corpus_words (integer): number of distinct words across the␣
      ↪corpus
          """
          corpus_words = []
          num_corpus_words = -1

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

          corpus_words = sorted(list(set(word for string in corpus for word in␣
      ↪string)))
          num_corpus_words = len(corpus_words)

          # ------------------

          return corpus_words, num_corpus_words
```

```
[37]: # ----------------------
      # Run this sanity check
      # Note that this not an exhaustive check for correctness.
      # ----------------------

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

```
# Correct answers
ans_test_corpus_words = sorted(list(set(["START", "All", "ends", "that",␣
 ↪"gold", "All's", "glitters", "isn't", "well", "END"])))
ans_num_corpus_words = len(ans_test_corpus_words)

# Test correct number of words
assert(num_corpus_words == ans_num_corpus_words), "Incorrect number of distinct␣
 ↪words. Correct: {}. Yours: {}".format(ans_num_corpus_words, num_corpus_words)

# Test correct words
assert (test_corpus_words == ans_test_corpus_words), "Incorrect corpus_words.
 ↪\nCorrect: {}\nYours:   {}".format(str(ans_test_corpus_words),␣
 ↪str(test_corpus_words))

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

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

### 2.1.4 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.

```
[38]: def compute_co_occurrence_matrix(corpus, window_size=4):
          """ Compute co-occurrence matrix for the given corpus and window_size␣
      ↪(default of 4).

              Note: Each word in a document should be at the center of a window.␣
      ↪Words near edges will have a smaller
                  number of co-occurring words.

                  For example, if we take the document "START All that glitters is␣
      ↪not gold END" with window size of 4,
                  "All" will co-occur with "START", "that", "glitters", "is", and␣
      ↪"not".

              Params:
                  corpus (list of list of strings): corpus of documents
                  window_size (int): size of context window
```

```
        Return:
            M (numpy matrix of shape (number of corpus words, number of corpus
    ↪words)):
                Co-occurence matrix of word counts.
                The ordering of the words in the rows/columns should be the
    ↪same as the ordering of the words given by the distinct_words function.
            word2Ind (dict): dictionary that maps word to index (i.e. row/
    ↪column number) for matrix M.
        """
    words, num_words = distinct_words(corpus)
    M = None
    word2Ind = {}

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

    M = np.zeros((num_words, num_words))
    word2Ind = {word:idx for idx, word in enumerate(words)}

    for string in corpus:
        for word_position in range(0, len(string)):
            for n in range(1, window_size + 1):
                if word_position-n >= 0:
                    M[word2Ind[string[word_position]],
    ↪word2Ind[string[word_position-n]]] += 1
                if word_position+n < len(string):
                    M[word2Ind[string[word_position]],
    ↪word2Ind[string[word_position+n]]] += 1

    # ------------------

    return M, word2Ind
```

```
[39]:   # ---------------------
        # Run this sanity check
        # Note that this is not an exhaustive check for correctness.
        # ---------------------

        # Define toy corpus and get student's co-occurrence matrix
        test_corpus = ["START All that glitters isn't gold END".split(" "), "START
        ↪All's well that ends well END".split(" ")]
        M_test, word2Ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)

        # Correct M and word2Ind
        M_test_ans = np.array(
            [[0., 0., 0., 1., 0., 0., 0., 0., 1., 0.,],
```

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

# Test correct word2Ind
assert (word2Ind_ans == word2Ind_test), "Your word2Ind is incorrect:\nCorrect:␣
 ↪{}\nYours: {}".format(word2Ind_ans, word2Ind_test)

# Test correct M shape
assert (M_test.shape == M_test_ans.shape), "M matrix has incorrect shape.
 ↪\nCorrect: {}\nYours: {}".format(M_test.shape, M_test_ans.shape)

# Test correct M values
for w1 in word2Ind_ans.keys():
    idx1 = word2Ind_ans[w1]
    for w2 in word2Ind_ans.keys():
        idx2 = word2Ind_ans[w2]
        student = M_test[idx1, idx2]
        correct = M_test_ans[idx1, idx2]
        if student != correct:
            print("Correct M:")
            print(M_test_ans)
            print("Your M: ")
            print(M_test)
            raise AssertionError("Incorrect count at index ({}, {})=({}, {}) in␣
 ↪matrix M. Yours has {} but should have {}.".format(idx1, idx2, w1, w2,␣
 ↪student, correct))

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

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

### 2.1.5 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.

```python
[40]: def reduce_to_k_dim(M, k=2):
          """ Reduce a co-occurence count matrix of dimensionality (num_corpus_words,␣
      ↪num_corpus_words)
              to a matrix of dimensionality (num_corpus_words, k) using the following␣
      ↪SVD function from Scikit-Learn:
                  - http://scikit-learn.org/stable/modules/generated/sklearn.
      ↪decomposition.TruncatedSVD.html

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

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

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

          # ------------------

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

```python
[41]: # ---------------------
      # 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.
```

14

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

# Define toy corpus and run student code
test_corpus = ["START All that glitters isn't gold END".split(" "), "START␣
 ↪All's well that ends well END".split(" ")]
M_test, word2Ind_test = compute_co_occurrence_matrix(test_corpus, window_size=1)
M_test_reduced = reduce_to_k_dim(M_test, k=2)

# Test proper dimensions
assert (M_test_reduced.shape[0] == 10), "M_reduced has {} rows; should have {}".
 ↪format(M_test_reduced.shape[0], 10)
assert (M_test_reduced.shape[1] == 2), "M_reduced has {} columns; should have␣
 ↪{}".format(M_test_reduced.shape[1], 2)

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

```
Running Truncated SVD over 10 words…
Done.
--------------------------------------------------------------------------------
Passed All Tests!
--------------------------------------------------------------------------------
```

### 2.1.6  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.

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

              Params:
                  M_reduced (numpy matrix of shape (number of unique words in the␣
      ↪corpus , k)): matrix of k-dimensioal word embeddings
                  word2Ind (dict): dictionary that maps word to indices for matrix M
                  words (list of strings): words whose embeddings we want to visualize
          """
```

15

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

        for word in words:
            x = M_reduced[word2Ind[word]][0]
            y = M_reduced[word2Ind[word]][1]
            plt.scatter(x, y, marker='x', color='red')
            plt.text(x, y, word, fontsize=9)

        plt.show()

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

[43]:
```
# --------------------
# Run this sanity check
# Note that this not an exhaustive check for correctness.
# The plot produced should look like the "test solution plot" depicted below.
# --------------------

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

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

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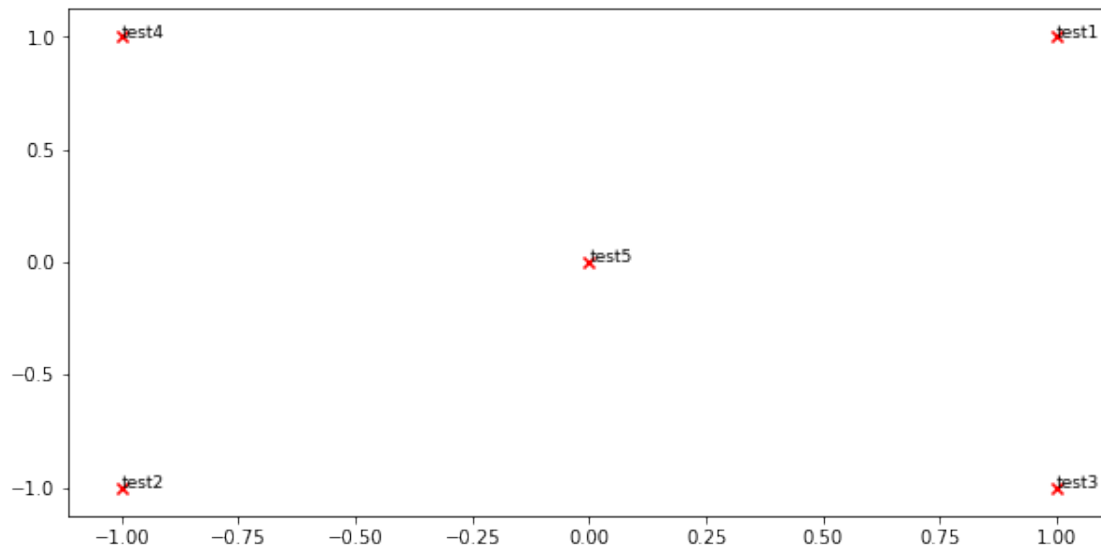
--------------------------------------------------------------------------------
Outputted Plot:

----------------------------------------------------------------------------------

**Test Plot Solution**

### 2.1.7 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 4, over the Reuters "crude" corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns U*S, so we normalize the returned vectors, so that all the vectors will appear around the unit circle (therefore closeness is directional closeness). **Note**: The line of code below that does the normalizing uses the NumPy concept of broadcasting. If you don't know about broadcasting, check out Computation on Arrays: Broadcasting by Jake VanderPlas.
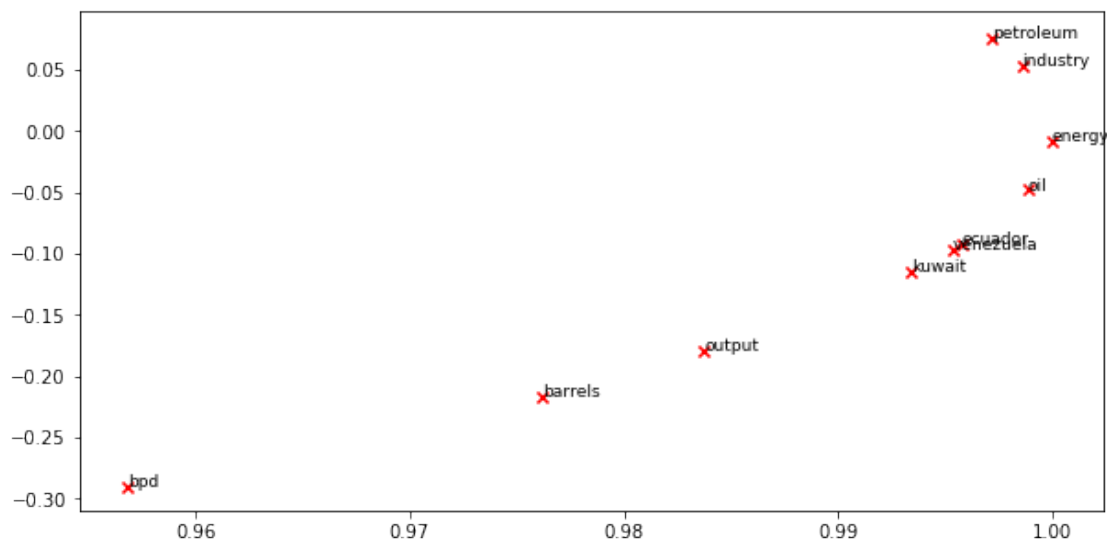
Run the below cell to produce the plot. It'll probably take a few seconds to run. What clusters together in 2-dimensional embedding space? What doesn't cluster together that you might think should have? **Note:** "bpd" stands for "barrels per day" and is a commonly used abbreviation in crude oil topic articles.

```
[44]:  # -----------------------------
       # 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
```

17

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

Running Truncated SVD over 8185 words…
Done.



1. What clusters together in 2-dimensional embedding space?
   - cluster1: ecuador, venezuela and kuwait
   - cluster2: petroleum and industry
   - cluster3: energy and oil
2. What doesn't cluster together that you might think should have?
   - cluster1: bpd, output and barrels

## 2.2 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.

```
[45]: def load_word2vec():
          """ Load Word2Vec Vectors
              Return:
                  wv_from_bin: All 3 million embeddings, each lengh 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
```

```
[46]: # --------------------------------
      # 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.**

### 2.2.1 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.

```
[47]: def get_matrix_of_vectors(wv_from_bin, required_words=['barrels', 'bpd',
      →'ecuador', 'energy', 'industry', 'kuwait', 'oil', 'output', 'petroleum',
      →'venezuela']):
          """ Put the word2vec vectors into a matrix M.
              Param:
                  wv_from_bin: KeyedVectors object; the 3 million word2vec vectors
      →loaded from file
              Return:
                  M: numpy matrix shape (num words, 300) containing the vectors
                  word2Ind: dictionary mapping each word to its row number in M
          """
          import random
          words = list(wv_from_bin.vocab.keys())
          print("Shuffling words ...")
          random.shuffle(words)
          words = words[:10000]
          print("Putting %i words into word2Ind and matrix M..." % len(words))
          word2Ind = {}
```

```
        M = []
        curInd = 0
        for w in words:
            try:
                M.append(wv_from_bin.word_vec(w))
                word2Ind[w] = curInd
                curInd += 1
            except KeyError:
                continue
        for w in required_words:
            try:
                M.append(wv_from_bin.word_vec(w))
                word2Ind[w] = curInd
                curInd += 1
            except KeyError:
                continue
        M = np.stack(M)
        print("Done.")
        return M, word2Ind
```

[48]:
```
# ----------------------------------------------------------------
# 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.
```

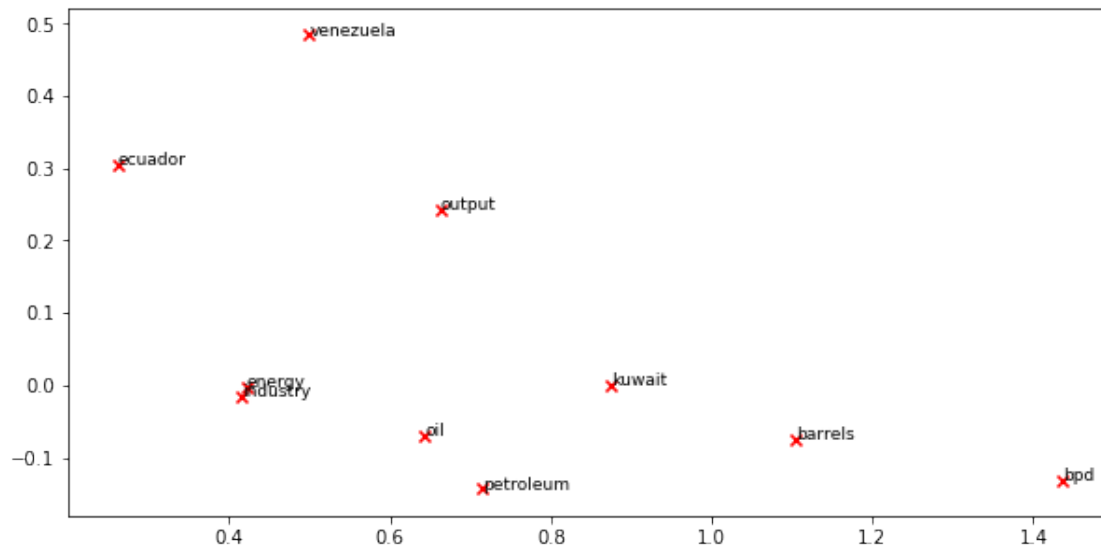### 2.2.2   Question 2.1: Word2Vec Plot Analysis [written] (4 points)

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 to-
gether that you might think should have? How is the plot different from the one generated
earlier from the co-occurrence matrix?

[49]:
```
words = ['barrels', 'bpd', 'ecuador', 'energy', 'industry', 'kuwait', 'oil',␣
 ↪'output', 'petroleum', 'venezuela']
plot_embeddings(M_reduced, word2Ind, words)
```

1. What clusters together in 2-dimensional embedding space?
   - energy and industry
2. What doesn't cluster together that you might think should have?
   - barrels, petroleum and bpd
   - ecuador, kuwait and venezuela
3. How is the plot different from the one generated earlier from the co-occurrence matrix?
   - Co-occurrence matrix' plot looks like a curve.
   - Word2Vec's plot is more sparse.
   - Word2Vec doesn't cluster words.

### 2.2.3  Cosine Similarity

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

We can think of n-dimensional vectors as points in n-dimensional space. If we take this perspective L1 and L2 Distances help quantify the amount of space "we must travel" to get between these two points. Another approach is to examine the angle between two vectors. From trigonometry we know that:

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

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

### 2.2.4 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**.

```
[50]: # ------------------
      # Write your polysemous word exploration code here.

      word = "cut"
      ans = wv_from_bin.most_similar(word)
      print(word+":")
      for topn in ans:
          print("\t"+str(topn))

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

```
cut:
        ('cutting', 0.7856094837188721)
        ('slash', 0.7417387962341309)
        ('slashed', 0.7039529085159302)
        ('trimmed', 0.6781732439994812)
        ('Cut', 0.6464383602142334)
        ('trimming', 0.6387491226196289)
        ('trim', 0.6330385208129883)
        ('slashing', 0.6282667517662048)
        ('cuts', 0.6245189309120178)
        ('lop', 0.598679780960083)
```

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

Word2Vec 的 embedding 表示該字在文集中出現的位置，和其他字的關係(co-occurrence)。透過 pre/post-fix 出現的頻率來表示這個字的意義，多義詞因為意思不同 pre/post-fix 出現的字可能不同。 若和其他同意的詞的 pre/post-fix 不相近的話，就會學不到。

舉例來說： "right" 的 top 10 中出現 "wrong" 和 "left"，不合格的多義詞 top 10。 因為 right 的 pre/post-fix 和 wrong/left 的類似，經過學習後，embedding 就會在向量空間中和這兩個字較接近。

### 2.2.5 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 - Cosine Similarity.

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

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

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

```
[51]:  # ------------------
       # Write your synonym & antonym exploration code here.


       w1 = "beautiful"
       w2 = "pretty"
       w3 = "ugly"
       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 beautiful, pretty have cosine distance: 0.6700939834117889
Antonyms beautiful, ugly have cosine distance: 0.6655565500259399
```

Why this counter-intuitive result may have happened?

有 "beautiful" 和 "ugly" 的句子中, "beautiful" 和 "ugly" 通常用來形容人物, 所以 pre/post-fix 出現的字可能比起 "beautiful" 和 "pretty" 的 pre/post-fix 較為相似。 "pretty" 除了形容人長得漂亮, 本身也能作為副詞 "非常" 的意思, 因此兩者常用的情境有更多的不同的可能, 但實際情況

### 2.2.6 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 **Gen-Sim documentation**.

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

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

### 2.2.7 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!

```
[53]: # ------------------
      # Write your analogy exploration code here.

      pprint.pprint(wv_from_bin.most_similar(positive=["Japan", "Taipei"],
       →negative=["Taiwan"]))

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

```
[('Tokyo', 0.7887160778045654),
 ('Nagoya', 0.6903098821640015),
 ('Yokohama', 0.6869075894355774),
 ('Osaka', 0.6853444576263428),
 ('Maebashi', 0.6708757877349854),
 ('Fukuoka', 0.6618070602416992),
 ('Saitama', 0.656663179397583),
 ('Chiba', 0.6446334719657898),
 ('Takamatsu', 0.6383306384086609),
 ('Tokyo_Chiyoda_Ward', 0.6247940063476562)]
```

Taiwan : Taipei :: Japan : Tokyo

The capital of Taiwan is Taipei and the capital of Japan is Tokyo.

### 2.2.8 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.

```
[54]:  # ------------------
       # Write your incorrect analogy exploration code here.

       pprint.pprint(wv_from_bin.most_similar(positive=["New_Zealand","Beijing"],
         →negative=["China"]))

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

```
[('Auckland', 0.720319390296936),
 ('Christchurch', 0.7060885429382324),
 ('NZ', 0.6936110258102417),
 ('Wellington', 0.6511342525482178),
 ('Kiwi', 0.6478166580200195),
 ('Canberra', 0.6419323682785034),
 ('Invercargill', 0.6328139305114746),
 ('New_Zealanders', 0.6290636658668518),
 ('Palmerston_North', 0.6235657930374146),
 ('Rotorua', 0.6184718012809753)]
```

China : Beijing :: New_Zealand : Auckland

Wellington has been the capital of New Zealand since 1865.

### 2.2.9 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?

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

```
[('bosses', 0.5522644519805908),
 ('manageress', 0.49151360988616943),
```

```
('exec', 0.45940813422203064),
('Manageress', 0.45598435401916504),
('receptionist', 0.4474116563796997),
('Jane_Danson', 0.44480544328689575),
('Fiz_Jennie_McAlpine', 0.44275766611099243),
('Coronation_Street_actress', 0.44275566935539246),
('supremo', 0.4409853219985962),
('coworker', 0.43986251950263977)]

[('supremo', 0.6097398400306702),
('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', 0.5256162881851196),
('striker_Gabby_Agbonlahor', 0.5250812768936157),
('BARNSLEY_boss', 0.5238943099975586)]
```

man : boss :: woman : x

x 可分為幾類: * 執行: exec、coworker * 接待員: receptionist * Coronation Street 內的女演員: Jane_Danson、Fiz_Jennie_McAlpine、Coronation_Street_actress * 接近 boss 的意思: supremo、Manageress、bosses

woman : boss :: man : y

y 可分為幾類: * 接近 boss 的意思: supremo、head_honcho * 球隊老闆: MOTHERWELL_boss、CARETAKER_boss、Bully_Wee_boss、YEOVIL_Town_boss、manager_Sta

corpus 明顯有性別偏差。男人是 boss 對比女人，確出現執行者、接待員、演過老闆的女演員；反過來女人是 boss 對比男人，就出現了一堆球隊的經理，和負責人等角色。

### 2.2.10 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.

```python
[56]: # ------------------
      # Write your bias exploration code here.

      pprint.pprint(wv_from_bin.most_similar(positive=["woman","Housework"],␣
       →negative=["man"]))
      print()
      pprint.pprint(wv_from_bin.most_similar(positive=["man","Housework"],␣
       →negative=["woman"]))

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

```
[('housework', 0.5822943449020386),
 ('Sex_Habits', 0.5090413093566895),
 ('ITV1_Loose', 0.5005926489830017),
 ('motherhood', 0.4949650168418884),
 ('Lactating', 0.4922730028629303),
 ('Motherhood', 0.49087172746658325),
 ('Yummy_Mummy', 0.4897937774658203),
 ('child_rearing', 0.4896451234817505),
 ('Unwed', 0.48835599422454834),
 ('Breastfeed', 0.486355185508728)]

[('THE_BRAZEN_CAREERIST', 0.47392332553863525),
 ('Whining', 0.45564693212509155),
 ('Sam_Seboe_column', 0.45538777112960815),
 ('Chores', 0.4479844868183136),
 ('Idleness', 0.44395607709884644),
 ('Cinematical_Seven', 0.4419426918029785),
 ('STAWAR', 0.441253662109375),
 ('Chianca', 0.44096094369888306),
 ('Boredom', 0.43809986114501953),
 ('Self_deprecating_humor', 0.4372100234031677)]
```

man : Housework :: woman : x

x 可分為幾類: * 育兒: motherhood、Lactating、Motherhood、child_rearing、Breastfeed * 電視節目: ITV1_Loose、Yummy_Mummy * 家事: housework * 個人: Unwed、Sex_Habits

woman : Housework :: man : y

y 可分為幾類: * 無關的書籍: THE_BRAZEN_CAREERIST、Cinematical_Seven * 情緒: Whining、Idleness、Boredom、Self_deprecating_humor * 雜事: Chores * 玩樂: STAWAR、Chianca

corpus 明顯有性別偏差。男人是 Housework 對比女人，養育孩子的責任、家事等等；反過來女人是 Housework 對比男人，就出現了一堆玩樂，抱怨的情緒，甚至是星戰。

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

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

Training data 的問題，所有的機器學習，都是基於餵給機器的 input 去學習，並產生 model 來描述 input，透過 model 產生的 output 就會符合 input 的資料。因為 demo 的資料集是基於 google news，觀察 Question 2.6 和 2.7，發現很多性別刻板印象偏差，我想多少反映了新聞上，或是普遍社會對男女性別的錯誤期待。

# 3  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 ll 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.