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

## Due 3:15pm, Tue Jan 11

**Welcome to CS224N!**

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 review session on Friday. The session will be recorded and the material will be made available on our website. The CS231N Python/Numpy tutorial is also a great resource.

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

In [1]:

```python
# All Import Statements Defined Here
# Note: Do not add to this list.
# ---------------

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

from platform import python_version
assert int(python_version().split(".")[1]) >= 5, "Please upgrade your Python version following the instructions in \
    the README.txt file found in the same directory as this notebook. Your Python version is " + python_version()

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') #to specify download location, optionally add the argument: download_dir='/specify/desired/pa
th/'
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 /home/mantas/nltk_data...
```

# 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 *GloVe*.

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

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

Most word vector models start from the following idea:

*You shall know a word by the company it keeps (Firth, J. R. 1957:11)*

Many word vector implementations are driven by the idea that similar words, i.e., (near) synonyms, will be used in similar contexts. As a result, similar words will often be spoken or written along with a shared subset of words, i.e., contexts. By examining these contexts, we can try to develop embeddings for our words. With this intuition in mind, many "old school" approaches to constructing word vectors relied on word counts. Here we elaborate upon one of those strategies, *co-occurrence matrices* (for more information, see here or here).

## Co-Occurrence

A co-occurrence matrix counts how often things co-occur in some environment. Given some word $w_i$ occurring in the document, we consider the *context window* surrounding $w_i$. Supposing our fixed window size is $n$, then this is the $n$ preceding and $n$ subsequent words in that document, i.e. words $w_{i-n} \ldots w_{i-1}$ 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 among all documents.

**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. If you want to learn more thoroughly about PCA or SVD, feel free to check out lectures 7, 8, and 9 of CS168. These course notes provide a great high-level treatment of these general purpose algorithms. Though, for the purpose of this class, you only need to know how to extract the k-dimensional embeddings by utilizing pre-programmed implementations of these algorithms from the numpy, scipy, or sklearn python packages. In practice, it is challenging to apply full SVD to large corpora because of the memory needed to perform PCA or SVD. However, if you only want the top $k$ vector components for relatively small $k$ — known as **Truncated SVD** — then there are reasonably scalable techniques to compute those iteratively.

## Plotting Co-Occurrence Word Embeddings

Here, we will be using the Reuters (business and financial news) corpus. If you haven't run the import cell at the top of this page, please run it now (click it and press 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 "grain" (i.e. news articles about corn, wheat, etc.) category. The function also adds <START> and <END> tokens to each of the documents, and lowercases words. You do **not** have to perform any other kind of pre-processing.

In [2]:

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

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

In [3]:

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

```
[['<START>', 'china', 'daily', 'says', 'vermin', 'eat', '7', '-', '12', 'pct', 'grain', 'stocks',
  'a', 'survey', 'of', '19', 'provinces', 'and', 'seven', 'cities', 'showed', 'vermin', 'consume',
  'between', 'seven', 'and', '12', 'pct', 'of', 'china', '"', 's', 'grain', 'stocks', ',', 'the',
  'china', 'daily', 'said', '.', 'it', 'also', 'said', 'that', 'each', 'year', '1', '.', '575',
  'mln', 'tonnes', ',', 'or', '25', 'pct', ',', 'of', 'china', '"', 's', 'fruit', 'output', 'are',
  'left', 'to', 'rot', ',', 'and', '2', '.', '1', 'mln', 'tonnes', ',', 'or', 'up', 'to', '30',
  'pct', ',', 'of', 'its', 'vegetables', '.', 'the', 'paper', 'blamed', 'the', 'waste', 'on',
  'inadequate', 'storage', 'and', 'bad', 'preservation', 'methods', '.', 'it', 'said', 'the',
  'government', 'had', 'launched', 'a', 'national', 'programme', 'to', 'reduce', 'waste', ',',
  'calling', 'for', 'improved', 'technology', 'in', 'storage', 'and', 'preservation', ',', 'and',
  'greater', 'production', 'of', 'additives', '.', 'the', 'paper', 'gave', 'no', 'further',
  'details', '.', '<END>'],
 ['<START>', 'thai', 'trade', 'deficit', 'widens', 'in', 'first', 'quarter', 'thailand', '"', 's',
  'trade', 'deficit', 'widened', 'to', '4', '.', '5', 'billion', 'baht', 'in', 'the', 'first',
  'quarter', 'of', '1987', 'from', '2', '.', '1', 'billion', 'a', 'year', 'ago', ',', 'the',
  'business', 'economics', 'department', 'said', '.', 'it', 'said', 'janunary', '/', 'march',
  'imports', 'rose', 'to', '65', '.', '1', 'billion', 'baht', 'from', '58', '.', '7', 'billion',
  '.', 'thailand', '"', 's', 'improved', 'business', 'climate', 'this', 'year', 'resulted', 'in',
  'a', '27', 'pct', 'increase', 'in', 'imports', 'of', 'raw', 'materials', 'and', 'semi', '-',
  'finished', 'products', '.', 'the', 'country', '"', 's', 'oil', 'import', 'bill', ',', 'however',
  '.', 'fell', '23', 'pct', 'in', 'the', 'first', 'quarter', 'due', 'to', 'lower', 'oil', 'prices'
```

```
 ',', 'tell', '25', 'pct', 'in', 'the', 'first', 'quarter', 'due', 'to', 'lower', 'oil', 'prices',
 '.', 'the', 'department', 'said', 'first', 'quarter', 'exports', 'expanded', 'to', '60', '.', '6',
 'billion', 'baht', 'from', '56', '.', '6', 'billion', '.', 'export', 'growth', 'was', 'smaller',
 'than', 'expected', 'due', 'to', 'lower', 'earnings', 'from', 'many', 'key', 'commodities',
 'including', 'rice', 'whose', 'earnings', 'declined', '18', 'pct', ',', 'maize', '66', 'pct', ',',
 'sugar', '45', 'pct', ',', 'tin', '26', 'pct', 'and', 'canned', 'pineapples', 'seven', 'pct', '.',
 'products', 'registering', 'high', 'export', 'growth', 'were', 'jewellery', 'up', '64', 'pct',
 ',', 'clothing', '57', 'pct', 'and', 'rubber', '35', 'pct', '.', '<END>'],
['<START>', 'sri', 'lanka', 'gets', 'usda', 'approval', 'for', 'wheat', 'price', 'food',
 'department', 'officials', 'said', 'the', 'u', '.', 's', '.', 'department', 'of', 'agriculture',
 'approved', 'the', 'continental', 'grain', 'co', 'sale', 'of', '52', ',', '500', 'tonnes', 'of',
 'soft', 'wheat', 'at', '89', 'u', '.', 's', '.', 'dlrs', 'a', 'tonne', 'c', 'and', 'f', 'from',
 'pacific', 'northwest', 'to', 'colombo', '.', 'they', 'said', 'the', 'shipment', 'was', 'for',
 'april', '8', 'to', '20', 'delivery', '.', '<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 may be useful to flatten a list of lists. If you're not familiar with Python list comprehensions in general, here's more information.

Your returned `corpus_words` should be sorted. You can use python's `sorted` function for this.

You may find it useful to use Python sets to remove duplicate words.

In [4]:

```python
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): sorted list of distinct words across the corpus
            n_corpus_words (integer): number of distinct words across the corpus
    """
    corpus_words = []
    n_corpus_words = -1

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

    corpus_words = sorted(set(sum(corpus, [])))
    n_corpus_words = len(corpus_words)

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

    return corpus_words, n_corpus_words
```

In [5]:

```python
# ---------------------
# Run this sanity check
# Note that this not an exhaustive check for correctness.
# ---------------------

# Define toy corpus
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that ends well {}".format(START_TOKEN, END_TOKEN).split(" ")]
test_corpus_words, num_corpus_words = distinct_words(test_corpus)

# Correct answers
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_TOKEN])
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!
--------------------------------------------------------------------------------
```

## 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.

In [6]:

```python
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 (a symmetric numpy matrix of shape (number of unique words in the corpus , number of unique words in the
corpus)):
                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, n_words = distinct_words(corpus)
    M = None
    word2ind = {}

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

    # Word to index mapping and a marix with 0s
    word2ind = dict(zip(words, range(n_words)))
    M = np.zeros((n_words, n_words))

    for document in corpus:
        for i, center in enumerate(document):
            for j, context in enumerate(document):
                if j != i and abs(j - i) <= window_size:
                    # Increment the co-occurance of word pair
                    M[word2ind[center]][word2ind[context]] += 1

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

    return M, word2ind
```

In [7]:

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

# Define toy corpus and get student's co-occurrence matrix
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that end
s well {}".format(START_TOKEN, END_TOKEN).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., 0., 0., 0., 1., 0., 0., 1.,],
     [0., 0., 1., 1., 0., 0., 0., 0., 0., 0.,],
     [0., 1., 0., 0., 0., 0., 0., 0., 1., 0.,],
     [0., 1., 0., 0., 0., 0., 0., 0., 0., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 0., 1., 1.,],
     [0., 0., 0., 0., 0., 0., 0., 1., 1., 0.,],
     [1., 0., 0., 0., 0., 0., 0., 1., 0., 0.,],
     [0., 0., 0., 0., 0., 1., 1., 0., 0., 0.,],
     [0., 0., 1., 0., 1., 1., 0., 0., 0., 1.,],
     [1., 0., 0., 1., 1., 0., 0., 0., 1., 0.,]]
)
ans_test_corpus_words = sorted([START_TOKEN, "All", "ends", "that", "gold", "All's", "glitters", "isn't", "well", END_
TOKEN])
word2ind_ans = dict(zip(ans_test_corpus_words, range(len(ans_test_corpus_words))))

# Test correct word2ind
assert (word2ind_ans == word2ind_test), "Your word2ind is incorrect:\nCorrect: {}\nYours: {}".format(word2ind_ans, wor
d2ind_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!
```

---

## 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.

In [8]:

```python
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 unique words in the corpus , number of unique words in the corpus)): c
o-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.

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

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

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

In [9]:

```python
# --------------------
# Run this sanity check
# Note that this is not an exhaustive check for correctness
# In fact we only check that your M_reduced has the right dimensions.
# --------------------

# Define toy corpus and run student code
test_corpus = ["{} All that glitters isn't gold {}".format(START_TOKEN, END_TOKEN).split(" "), "{} All's well that end
s well {}".format(START_TOKEN, END_TOKEN).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!
--------------------------------------------------------------------------------
```

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

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

For this example, you may find it useful to adapt this code. In the future, a good way to make a plot is to look at the Matplotlib gallery, find a plot that looks somewhat like what you want, and adapt the code they give.

In [10]:

```python
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 , 2)): matrix of 2-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
    """

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

```
    # Use a prettier seaborn grid style
    plt.style.use('seaborn-whitegrid')

    for i, word in enumerate(words):
        # Get the coordinates of the embedding
        [x, y] = M_reduced[word2ind[word], :]

        # Mark the coordinate and assign text to it
        plt.scatter(x, y, marker='x', color="red")
        plt.annotate(word, (x, y), xytext=(x, y+.05))

    # Show plot
    plt.show()

    # ------------------
```
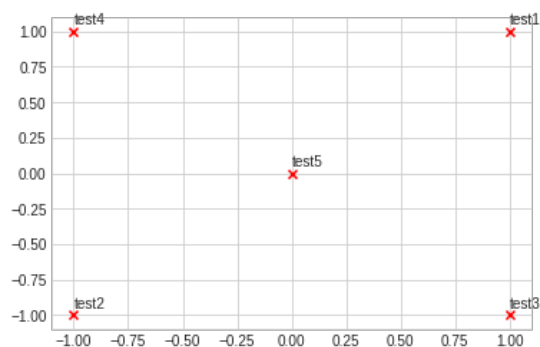
In [11]:

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

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

M_reduced_plot_test = np.array([[1, 1], [-1, -1], [1, -1], [-1, 1], [0, 0]])
word2ind_plot_test = {'test1': 0, 'test2': 1, 'test3': 2, 'test4': 3, 'test5': 4}
words = ['test1', 'test2', 'test3', 'test4', 'test5']
plot_embeddings(M_reduced_plot_test, word2ind_plot_test, words)

print ("-" * 80)
```

--------------------------------------------------------------------------------
Outputted Plot:



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

**Test Plot Solution**

## Question 1.5: Co-Occurrence Plot Analysis [written] (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 (the default window size), over the Reuters "grain" corpus. Then we will use TruncatedSVD to compute 2-dimensional embeddings of each word. TruncatedSVD returns U*S, so we need to 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?

In [12]:

```
# ------------------------------
# Run This Cell to Produce Your Plot
# ------------------------------
reuters_corpus = read_corpus()
M_co_occurrence, word2ind_co_occurrence = compute_co_occurrence_matrix(reuters_corpus)
M_reduced_co_occurrence = reduce_to_k_dim(M_co_occurrence, k=2)

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

words = ['tonnes', 'grain', 'wheat',  'agriculture', 'corn', 'maize', 'export', 'department', 'barley', 'grains', 'soy
beans', 'sorghum']

plot_embeddings(M_normalized, word2ind_co_occurrence, words)
```
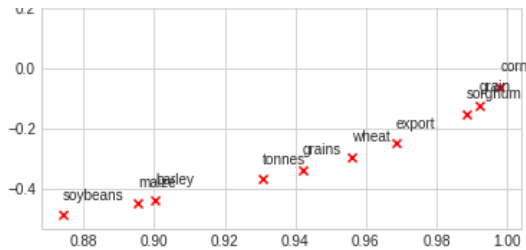
Running Truncated SVD over 7146 words...
Done.

1. Considering the scale of $x$, all the words seem close together which is reasonable as they are mainly the names of some sort of grain. Considering the scale of $y$ one cluster is [ *agriculture*, *department*] while the other consists of the remaining words which is also quite reasonable since the first cluster doesn't contain grain names. It should be noted, however, that the grain names are spread more in $y$ axis $(-0.5, 0)$ than in $x$ axis $(0.87, 1)$ which may seem counter-intuitive when looking at the plot.
2. If *grain* and *grains* were closer, it would be more logical as one is just a plural form of another word. It would also seem that *tonnes* and *export* should be closer to *department* as they are not the names of any kind of grain. Perhaps *2D* space is not sufficient to show they are similar.

## Part 2: Prediction-Based Word Vectors (15 points)

As discussed in class, more recently prediction-based word vectors have demonstrated better performance, such as word2vec and GloVe (which also utilizes the benefit of counts). Here, we shall explore the embeddings produced by GloVe. Please revisit the class notes and lecture slides for more details on the word2vec and GloVe algorithms. If you're feeling adventurous, challenge yourself and try reading GloVe's original paper.

Then run the following cells to load the GloVe vectors into memory. **Note:** If this is your first time to run these cells, i.e. download the embedding model, it will take a couple minutes to run. If you've run these cells before, rerunning them will load the model without redownloading it, which will take about 1 to 2 minutes.

In [13]:

```
def load_embedding_model():
    """ Load GloVe Vectors
        Return:
            wv_from_bin: All 400000 embeddings, each lengh 200
    """
    import gensim.downloader as api
    wv_from_bin = api.load("glove-wiki-gigaword-200")
    print("Loaded vocab size %i" % len(list(wv_from_bin.index_to_key)))
    return wv_from_bin
```

In [14]:

```
# -----------------------------------
# Run Cell to Load Word Vectors
# Note: This will take a couple minutes
# -----------------------------------
wv_from_bin = load_embedding_model()
```

[=================================================-] 98.2% 247.5/252.1MB downloadedLoaded vocab size 400000

**Note: If you are receiving a "reset by peer" error, rerun the cell to restart the download.**

### Reducing dimensionality of Word Embeddings

Let's directly compare the GloVe embeddings to those of the co-occurrence matrix. In order to avoid running out of memory, we will work with a sample of 10000 GloVe vectors instead. Run the following cells to:

1. Put 10000 Glove vectors into a matrix M
2. Run `reduce_to_k_dim` (your Truncated SVD function) to reduce the vectors from 200-dimensional to 2-dimensional.

In [15]:

```
def get_matrix_of_vectors(wv_from_bin, required_words=['tonnes', 'grain', 'wheat', 'agriculture', 'corn', 'maize', 'export', 'department', 'barley', 'grains', 'soybeans', 'sorghum']):
    """ Put the GloVe vectors into a matrix M.
        Param:
            wv_from_bin: KeyedVectors object; the 400000 GloVe vectors loaded from file
        Return:
            M: numpy matrix shape (num words, 200) containing the vectors
            word2ind: dictionary mapping each word to its row number in M
    """
    import random
    words = list(wv_from_bin.index_to_key)
    print("Shuffling words ...")
    random.seed(225)
    random.shuffle(words)
    words = words[:10000]
    print("Putting %i words into word2ind and matrix M..." % len(words))
    word2ind = {}
    M = []
    curInd = 0
    for w in words:
        try:
            M.append(wv_from_bin.get_vector(w))
            word2ind[w] = curInd
```

```
                curInd += 1
        except KeyError:
            continue
    for w in required_words:
        if w in words:
            continue
        try:
            M.append(wv_from_bin.get_vector(w))
            word2ind[w] = curInd
            curInd += 1
        except KeyError:
            continue
    M = np.stack(M)
    print("Done.")
    return M, word2ind
```

In [16]:

```python
# -----------------------------------------------------------------
# Run Cell to Reduce 200-Dimensional Word Embeddings to k Dimensions
# Note: This should be quick to run
# -----------------------------------------------------------------
M, word2ind = get_matrix_of_vectors(wv_from_bin)
M_reduced = reduce_to_k_dim(M, k=2)

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

```
Shuffling words ...
Putting 10000 words into word2ind and matrix M...
Done.
Running Truncated SVD over 10012 words...
Done.
```

**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 go to office hours or contact course staff.**

## Question 2.1: GloVe Plot Analysis [written] (3 points)

Run the cell below to plot the 2D GloVe embeddings for `['tonnes', 'grain', 'wheat', 'agriculture', 'corn', 'maize', 'export', 'department', 'barley', 'grains', 'soybeans', 'sorghum']`.
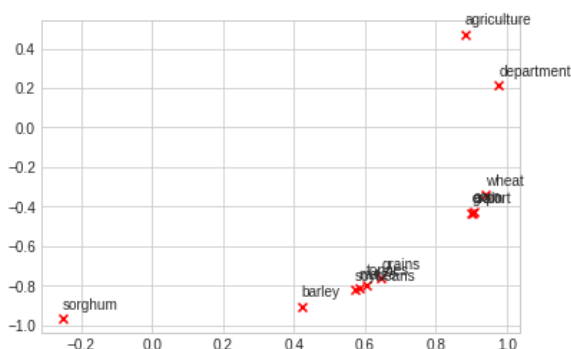
What clusters together in 2-dimensional embedding space? What doesn't cluster together that you think should have? How is the plot different from the one generated earlier from the co-occurrence matrix? What is a possible cause for the difference?

In [17]:

```python
words = ['tonnes', 'grain', 'wheat', 'agriculture', 'corn', 'maize', 'export', 'department', 'barley', 'grains', 'soybeans', 'sorghum']
plot_embeddings(M_reduced_normalized, word2ind, words)
```



**Write your answer here.**

1. There are 2 relatively dense clusters [ *grains*, *tonnes*, *soybeans*, *maize* ] and [ *wheat*, *grain*, *corn*, *export* ]. It is reasonable because both clusters capture similar category of words (excluding "*export*"), just that the words are not used as common in the first cluster as in the second.
2. Perhaps again *department*, *export* and *tonnes* should cluster together as they are not the names of any kind of grain. And *grains* and *grain* should be closer together as they are the same words just of the different forms. Also *sorghum* could also be closer to other kinds of grain.
3. The $x$ axis extends all the way to $-0.3$ (compared to $+0.87$ in the previous plot). The $y$ axis also has a wider range which starts at $-1$ (comared to $-0.5$ in the previous plot). Scatter-wise the shape of the plot is similar, just that $x$ coordinate is scaled to the negative side. There are some minor variations between word coordinates, e.g., *sorghum* is not close to *corn*.
4. **GloVe** generates word embeddings based on wikipedia which is used for training. It is more general which explains scale differences. Also, words like *sorghum* and *barley* are used there less frequently than in the news articles of "grain" category which explains why the similarities to other grain types are captured less accurately.

## 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:

$$= cos(\Theta)$$

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

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

## Question 2.2: Words with Multiple Meanings (1.5 points) [code + written]

Polysemes and homonyms are words that have more than one meaning (see this  wiki page  to learn more about the difference between polysemes and homonyms ). Find a word with *at least two different meanings* such that the top-10 most similar words (according to cosine similarity) contain related words from *both* meanings. For example, "leaves" has both "go_away" and "a_structure_of_a_plant" meaning in the top 10, and "scoop" has both "handed_waffle_cone" and "lowdown". You will probably need to try several polysemous or homonymic words before you find one.

Please state the word you discover and the multiple meanings that occur in the top 10. Why do you think many of the polysemous or homonymic words you tried didn't work (i.e. the top-10 most similar words only contain **one** of the meanings of the words)?

**Note:** You should use the `wv_from_bin.most_similar(word)` function to get the top 10 similar words. This function ranks all other words in the vocabulary with respect to their cosine similarity to the given word. For further assistance, please check the GenSim documentation.

In [18]:

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

wv_from_bin.most_similar("close")

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

Out[18]:

```
[('up', 0.6898430585861206),
 ('closed', 0.6854838728904724),
 ('down', 0.6749922633171082),
 ('while', 0.6741015911102295),
 ('just', 0.6566178798675537),
 ('but', 0.6522998213768005),
 ('closing', 0.6294093728065491),
 ('point', 0.6282009482383728),
 ('far', 0.6234282851219177),
 ('time', 0.6206085085868835)]
```

**Write your answer here.**

Homonym: **close** (*near/ not open*). Related words:

| Meaning 1 | Meaning 2 |
| --- | --- |
| Up | Closed |
| While | Down |
| But | Just |
| Point | Closing |
| Far | Time |

Some homonyms didn't work because one of the meanings is used more frequently than another. For instance, the word  **current** is used more often in the meaning of *present*, *up to date* than its noun version *stream*, as in *water current* or *electrical current*. Frequency also depends on the form of the word. For instance, **fly** is mainly associated with *flight*, *floating in air* whereas **flies** is a more specific version of the same word and is not used as frequently as the general form which allows to associate **flies** with *insects*.

## 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 $(w_1, w_2, w_3)$ where $w_1$ and $w_2$ are synonyms and $w_1$ and $w_3$ are antonyms, but Cosine Distance $(w_1, w_3) <$ Cosine Distance $(w_1, w_2)$.

As an example, $w_1$="happy" is closer to $w_3$="sad" than to $w_2$="cheerful". Please find a different example that satisfies the above. Once you have found your example, please give a possible explanation for why this counter-intuitive result may have happened.

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

In [19]:

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

W1 = "close"
W2 = "near"
```

```
W3 = "far"

print("Distance to synonym:", wv_from_bin.distance(W1, W2))
print("Distance to antonym:", wv_from_bin.distance(W1, W3))

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

```
Distance to synonym: 0.4232409596443176
Distance to antonym: 0.3765716552734375
```

**Write your answer here.**

$w_1 = \text{close}$ **is closer to** $w_3 = \text{far}$ **than to** $w_2 = \text{near}$**. The reason for such occurance is that** *close* **is used more often than** *near* **(like in this answer). Which means corpus has more occurances of** *far* **to support the context (which is "distance") than occurances of** *near* **to support the similarity.**

## Question 2.4: Analogies with Word Vectors [written] (1.5 points)

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

As an example, for the analogy "man : grandfather :: woman : x" (read: man is to grandfather as woman is to x), what is x?

In the cell below, we show you how to use word vectors to find x using the `most_similar` function from the [GenSim documentation](). The function finds words that are most similar to the words in the `positive` list and most dissimilar from the words in the `negative` list (while omitting the input words, which are often the most similar; see [this paper]()). The answer to the analogy will have the highest cosine similarity (largest returned numerical value).

In [20]:

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

```
[('grandmother', 0.7608445286750793),
 ('granddaughter', 0.7200808525085449),
 ('daughter', 0.7168302536010742),
 ('mother', 0.7151536345481873),
 ('niece', 0.7005682587623596),
 ('father', 0.6659887433052063),
 ('aunt', 0.6623408794403076),
 ('grandson', 0.6618767976760864),
 ('grandparents', 0.644661009311676),
 ('wife', 0.6445354223251343)]
```

Let $m$, $g$, $w$, and $x$ denote the word vectors for `man`, `grandfather`, `woman`, and the answer, respectively. Using **only** vectors $m$, $g$, $w$, and the vector arithmetic operators $+$ and $-$ in your answer, what is the expression in which we are maximizing cosine similarity with $x$?

Hint: Recall that word vectors are simply multi-dimensional vectors that represent a word. It might help to draw out a 2D example using arbitrary locations of each vector. Where would `man` and `woman` lie in the coordinate plane relative to `grandfather` and the answer?

**Write your answer here.**

$$w + g - m$$

## Question 2.5: Finding Analogies [code + written] (1.5 points)

Find an example of analogy that holds according to these vectors (i.e. the intended word is ranked top). In your solution please state the full analogy in the form x:y :: a:b. If you believe the analogy is complicated, explain why the analogy holds in one or two sentences.

**Note:** You may have to try many analogies to find one that works!

In [21]:

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

word1 = "dark"
word2 = "light"
isto1 = "black"
isto2 = wv_from_bin.most_similar(positive=[isto1, word2], negative=[word1])[0][0]

print(f"{word1} is to {isto1} as {word2} is to {isto2}")

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

```
dark is to black as light is to white
```

**Write your answer here.**

```
dark : black :: light : white
```

**This holds because** *dark* **is associated with a dark colour, i.e.,** *black***, and** *light* **is associated with a bright colour, i.e.,** *white***.**

## Question 2.6: Incorrect Analogy [code + written] (1.5 points)

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

In [22]:

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

word1 = "night"
word2 = "moon"
isto1 = "day"
isto2 = wv_from_bin.most_similar(positive=[isto1, word2], negative=[word1])[0][0]

print(f"{word1} is to {isto1} as {word2} is to {isto2}")

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

```
night is to day as moon is to lunar
```

<span style="color:red">**Write your answer here.**</span>

**Intended form:** `night : moon :: day : sun`

**The obtained word** `b=lunar`

## Question 2.7: 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 in our word embeddings. Bias can be dangerous because it can reinforce stereotypes through applications that employ these models.

Run the cell below, to examine (a) which terms are most similar to "girl" and "toy" and most dissimilar to "boy", and (b) which terms are most similar to "boy" and "toy" and most dissimilar to "girl". Point out the difference between the list of female-associated words and the list of male-associated words, and explain how it is reflecting gender bias.

In [23]:

```
# 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=['girl', 'toy'], negative=['boy']))
print()
pprint.pprint(wv_from_bin.most_similar(positive=['boy', 'toy'], negative=['girl']))
```

```
[('toys', 0.7094953060150146),
 ('doll', 0.5932914614677429),
 ('dolls', 0.570662260055542),
 ('barbie', 0.5407706499099731),
 ('mattel', 0.5328551530838013),
 ('accessories', 0.5206909775733948),
 ('hasbro', 0.49227219820022583),
 ('jewelry', 0.47385692596435547),
 ('lego', 0.4690813422203064),
 ('apparel', 0.46136239171028137)]

[('toys', 0.71570885181427),
 ('hasbro', 0.5164632797241211),
 ('robot', 0.47317108511924744),
 ('pet', 0.4670490324497223),
 ('manufacturer', 0.4668163061141968),
 ('mattel', 0.4582391679286957),
 ('lego', 0.45811763405799866),
 ('miniature', 0.4441472291946411),
 ('makers', 0.44298243522644043),
 ('manufactured', 0.44275349378585815)]
```

<span style="color:red">**Write your answer here.**</span>

**Representative terms for girl (a) and boy (b):**

1. Doll, dolls, barbie, accessories, jewelry, apparel
2. Robot, makers

In general, female-generated words are more related to female gender whereas male-generaed words are more neutral. The existing differences in the generated words that are unique to each gender reflect the gender bias, for instance boy toys are more associated with *robots* and girl toys are more associated with *dolls*.

## Question 2.8: Independent Analysis of Bias in Word Vectors [code + written] (1 point)

Use the `most_similar` function to find another case where some bias is exhibited by the vectors. Please briefly explain the example of bias that you discover.

In [24]:

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

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

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

```
[('clothing', 0.6792929172515869),
 ('shoes', 0.5772418975830078),
 ('shirts', 0.5559267997741699),
 ('jeans', 0.5477117300033569),
```

```
('wearing', 0.5434848666191101),
('everything', 0.5434514880180359),
('uniforms', 0.5397155284881592),
('stuff', 0.5329745411872864),
('shirt', 0.5325403213500977),
('hats', 0.525122344493866)]

[('clothing', 0.7139159440994263),
('dress', 0.6586079001426697),
('dresses', 0.6504972577095032),
('underwear', 0.5851817727088928),
('wear', 0.5704255700111389),
('gowns', 0.5619208216667175),
('shoes', 0.561784565448761),
('accessories', 0.5580349564552307),
('fashions', 0.5478761792182922),
('garments', 0.5419012308120728)]
```

**Write your answer here.**

Gender bias is chosen again as it proves to work quite well (compared to, for example, age bias). This time the context is *clothes* instead of *toys*. And the results make sense because people associate female clothes with *dress*, *accessories* while male clothes more with *jeans*, *shirts*.

## Question 2.9: Thinking About Bias [written] (2 points)

Give one explanation of how bias gets into the word vectors. What is an experiment that you could do to test for or to measure this source of bias?

**Write your answer here.**

1. **GloVe** embeddings also incorporate word co-occurances meaning that given the same context but different subjects, different descriptive words may be used which would result in bias, i.e., descriptive words make the subject embedding meaningful.
2. A possible experiment is to create sentences with the same context but leave the subject blank. For example, the context could be "shopping for clothes" and then the types of clothes as descriptive words are mentioned, e.g., "suit", "tie". Based on that, different people could be asked to fill in the blank and we could see what kind of subject (e.g., "male"/"female") the majority of people are biased to. The results can be compared with word vector bias. The vector bias itself could be measured by averaging the co-occurance counts between a descriptive word (for instance, "tie") and the desired subjects (for instance, "woman"/"man") across random documents.

# Submission Instructions

1. Click the Save button at the top of the Jupyter Notebook.
2. Select Cell -> All Output -> Clear. This will clear all the outputs from all cells (but will keep the content of all cells).
3. Select Cell -> Run All. This will run all the cells in order, and will take several minutes.
4. Once you've rerun everything, select File -> Download as -> PDF via LaTeX (If you have trouble using "PDF via LaTex", you can also save the webpage as pdf. **Make sure all your solutions especially the coding parts are displayed in the pdf**, it's okay if the provided codes get cut off because lines are not wrapped in code cells).
5. Look at the PDF file and make sure all your solutions are there, displayed correctly. The PDF is the only thing your graders will see!
6. Submit your PDF on Gradescope.