# Deep Learning

## Assignment 5

The goal of this assignment is to train a Word2Vec skip-gram model over Text8 (http://mattmahoney.net/dc/textdata) data.

In [1]:
```python
# These are all the modules we'll be using later. Make sure you can import them
# before proceeding further.
%matplotlib inline
from __future__ import print_function
import collections
import math
import numpy as np
import os
import random
import tensorflow as tf
import zipfile
from matplotlib import pylab
from six.moves import range
from six.moves.urllib.request import urlretrieve
from sklearn.manifold import TSNE
```

Download the data from the source website if necessary.

In [2]:
```python
url = 'http://mattmahoney.net/dc/'

def maybe_download(filename, expected_bytes):
  """Download a file if not present, and make sure it's the right size."""
  if not os.path.exists(filename):
    filename, _ = urlretrieve(url + filename, filename)
  statinfo = os.stat(filename)
  if statinfo.st_size == expected_bytes:
    print('Found and verified %s' % filename)
  else:
    print(statinfo.st_size)
    raise Exception(
      'Failed to verify ' + filename + '. Can you get to it with a browser?')
  return filename

filename = maybe_download('text8.zip', 31344016)
```

```
Found and verified text8.zip
```

Read the data into a string.

In [3]:
```python
def read_data(filename):
  """Extract the first file enclosed in a zip file as a list of words"""
  with zipfile.ZipFile(filename) as f:
    data = tf.compat.as_str(f.read(f.namelist()[0])).split()
  return data

words = read_data(filename)
print('Data size %d' % len(words))
```

```
Data size 17005207
```

Build the dictionary and replace rare words with UNK token.

```
In [4]:  vocabulary_size = 50000

         def build_dataset(words):
           count = [['UNK', -1]]
           count.extend(collections.Counter(words).most_common(vocabulary_size - 1))
           dictionary = dict()
           for word, _ in count:
             dictionary[word] = len(dictionary)
           data = list()
           unk_count = 0
           for word in words:
             if word in dictionary:
               index = dictionary[word]
             else:
               index = 0  # dictionary['UNK']
               unk_count = unk_count + 1
             data.append(index)
           count[0][1] = unk_count
           reverse_dictionary = dict(zip(dictionary.values(), dictionary.keys()))
           return data, count, dictionary, reverse_dictionary

         data, count, dictionary, reverse_dictionary = build_dataset(words)
         print('Most common words (+UNK)', count[:5])
         print('Sample data', data[:10])
         del words  # Hint to reduce memory.
```

```
Most common words (+UNK) [['UNK', 418391], ('the', 1061396), ('of', 593677), ('and', 416629), ('one', 411764)]
Sample data [5239, 3084, 12, 6, 195, 2, 3137, 46, 59, 156]
```

Function to generate a training batch for the skip-gram model.

```
In [5]:  data_index = 0

         def generate_batch(batch_size, num_skips, skip_window):
           global data_index
           assert batch_size % num_skips == 0
           assert num_skips <= 2 * skip_window
           batch = np.ndarray(shape=(batch_size), dtype=np.int32)
           labels = np.ndarray(shape=(batch_size, 1), dtype=np.int32)
           span = 2 * skip_window + 1 # [ skip_window target skip_window ]
           buffer = collections.deque(maxlen=span)
           for _ in range(span):
             buffer.append(data[data_index])
             data_index = (data_index + 1) % len(data)
           for i in range(batch_size // num_skips):
             target = skip_window  # target label at the center of the buffer
             targets_to_avoid = [ skip_window ]
             for j in range(num_skips):
               while target in targets_to_avoid:
                 target = random.randint(0, span - 1)
               targets_to_avoid.append(target)
               batch[i * num_skips + j] = buffer[skip_window]
               labels[i * num_skips + j, 0] = buffer[target]
             buffer.append(data[data_index])
             data_index = (data_index + 1) % len(data)
           return batch, labels

         print('data:', [reverse_dictionary[di] for di in data[:8]])

         for num_skips, skip_window in [(2, 1), (4, 2)]:
             data_index = 0
             batch, labels = generate_batch(batch_size=8, num_skips=num_skips, skip_window=skip_window)
             print('\nwith num_skips = %d and skip_window = %d:' % (num_skips, skip_window))
             print('    batch:', [reverse_dictionary[bi] for bi in batch])
             print('    labels:', [reverse_dictionary[li] for li in labels.reshape(8)])
```

```
data: ['anarchism', 'originated', 'as', 'a', 'term', 'of', 'abuse', 'first']

with num_skips = 2 and skip_window = 1:
    batch: ['originated', 'originated', 'as', 'as', 'a', 'a', 'term', 'term']
    labels: ['anarchism', 'as', 'originated', 'a', 'as', 'term', 'a', 'of']
```

```
with num_skips = 4 and skip_window = 2:
    batch: ['as', 'as', 'as', 'as', 'a', 'a', 'a', 'a']
    labels: ['originated', 'anarchism', 'a', 'term', 'originated', 'term', 'of', 'as']
```

Train a skip-gram model.

```
In [6]: batch_size = 128
        embedding_size = 128 # Dimension of the embedding vector.
        skip_window = 1 # How many words to consider left and right.
        num_skips = 2 # How many times to reuse an input to generate a label.
        # We pick a random validation set to sample nearest neighbors. here we limit the
        # validation samples to the words that have a low numeric ID, which by
        # construction are also the most frequent.
        valid_size = 16 # Random set of words to evaluate similarity on.
        valid_window = 100 # Only pick dev samples in the head of the distribution.
        valid_examples = np.array(random.sample(range(valid_window), valid_size))
        num_sampled = 64 # Number of negative examples to sample.

        graph = tf.Graph()

        with graph.as_default(), tf.device('/cpu:0'):

          # Input data.
          train_dataset = tf.placeholder(tf.int32, shape=[batch_size])
          train_labels = tf.placeholder(tf.int32, shape=[batch_size, 1])
          valid_dataset = tf.constant(valid_examples, dtype=tf.int32)

          # Variables.
          embeddings = tf.Variable(
            tf.random_uniform([vocabulary_size, embedding_size], -1.0, 1.0))
          softmax_weights = tf.Variable(
            tf.truncated_normal([vocabulary_size, embedding_size],
                                stddev=1.0 / math.sqrt(embedding_size)))
          softmax_biases = tf.Variable(tf.zeros([vocabulary_size]))

          # Model.
          # Look up embeddings for inputs.
          embed = tf.nn.embedding_lookup(embeddings, train_dataset)
          # Compute the softmax loss, using a sample of the negative labels each time.
          loss = tf.reduce_mean(
            tf.nn.sampled_softmax_loss(weights=softmax_weights, biases=softmax_biases, inputs=embed,
                                       labels=train_labels, num_sampled=num_sampled, num_classes=vocabulary_size))

          # Optimizer.
          # Note: The optimizer will optimize the softmax_weights AND the embeddings.
          # This is because the embeddings are defined as a variable quantity and the
          # optimizer's `minimize` method will by default modify all variable quantities
```

```
# that contribute to the tensor it is passed.
# See docs on `tf.train.Optimizer.minimize()` for more details.
optimizer = tf.train.AdagradOptimizer(1.0).minimize(loss)

# Compute the similarity between minibatch examples and all embeddings.
# We use the cosine distance:
norm = tf.sqrt(tf.reduce_sum(tf.square(embeddings), 1, keep_dims=True))
normalized_embeddings = embeddings / norm
valid_embeddings = tf.nn.embedding_lookup(
  normalized_embeddings, valid_dataset)
similarity = tf.matmul(valid_embeddings, tf.transpose(normalized_embeddings))
```

In [7]:
```python
num_steps = 100001

with tf.Session(graph=graph) as session:
  tf.global_variables_initializer().run()
  print('Initialized')
  average_loss = 0
  for step in range(num_steps):
    batch_data, batch_labels = generate_batch(
      batch_size, num_skips, skip_window)
    feed_dict = {train_dataset : batch_data, train_labels : batch_labels}
    _, l = session.run([optimizer, loss], feed_dict=feed_dict)
    average_loss += l
    if step % 2000 == 0:
      if step > 0:
        average_loss = average_loss / 2000
      # The average loss is an estimate of the loss over the last 2000 batches.
      print('Average loss at step %d: %f' % (step, average_loss))
      average_loss = 0
    # note that this is expensive (~20% slowdown if computed every 500 steps)
    if step % 10000 == 0:
      sim = similarity.eval()
      for i in range(valid_size):
        valid_word = reverse_dictionary[valid_examples[i]]
        top_k = 8 # number of nearest neighbors
        nearest = (-sim[i, :]).argsort()[1:top_k+1]
        log = 'Nearest to %s:' % valid_word
        for k in range(top_k):
          close_word = reverse_dictionary[nearest[k]]
          log = '%s %s,' % (log, close_word)
        print(log)
  final_embeddings = normalized_embeddings.eval()
```

```
Nearest to at: during, near, licks, impregnated, around, reached, after, bayard,
Nearest to some: many, several, these, any, most, all, this, stained,
Nearest to were: are, had, was, have, those, while, these, homomorphisms,
Nearest to a: another, any, straps, lauren, sliced, elphinstone, every, the,
Nearest to to: would, will, should, integrating, cannot, must, hermaphroditus, surveyed,
Nearest to system: systems, process, foxes, theory, oppressive, devices, turnaround, concept,
Nearest to all: both, every, many, several, various, some, each, any,
Nearest to and: or, but, while, including, through, motorcycle, imagines, arr,
Nearest to one: seven, two, four, six, eight, five, three, nine,
Nearest to people: children, men, writers, women, students, person, casualties, jews,
```

```
Nearest to with: between, quintet, by, in, refitted, falkland, bowls, notebooks,
Nearest to its: their, his, the, her, your, adhesion, whose, ode,
Average loss at step 92000: 3.398948
Average loss at step 94000: 3.257952
Average loss at step 96000: 3.361469
Average loss at step 98000: 3.243324
Average loss at step 100000: 3.354439
Nearest to there: they, it, he, now, we, still, often, generally,
Nearest to d: b, j, elektra, morgan, convened, calvinists, archibald, giscard,
Nearest to have: had, has, having, be, are, tend, refer, were,
```
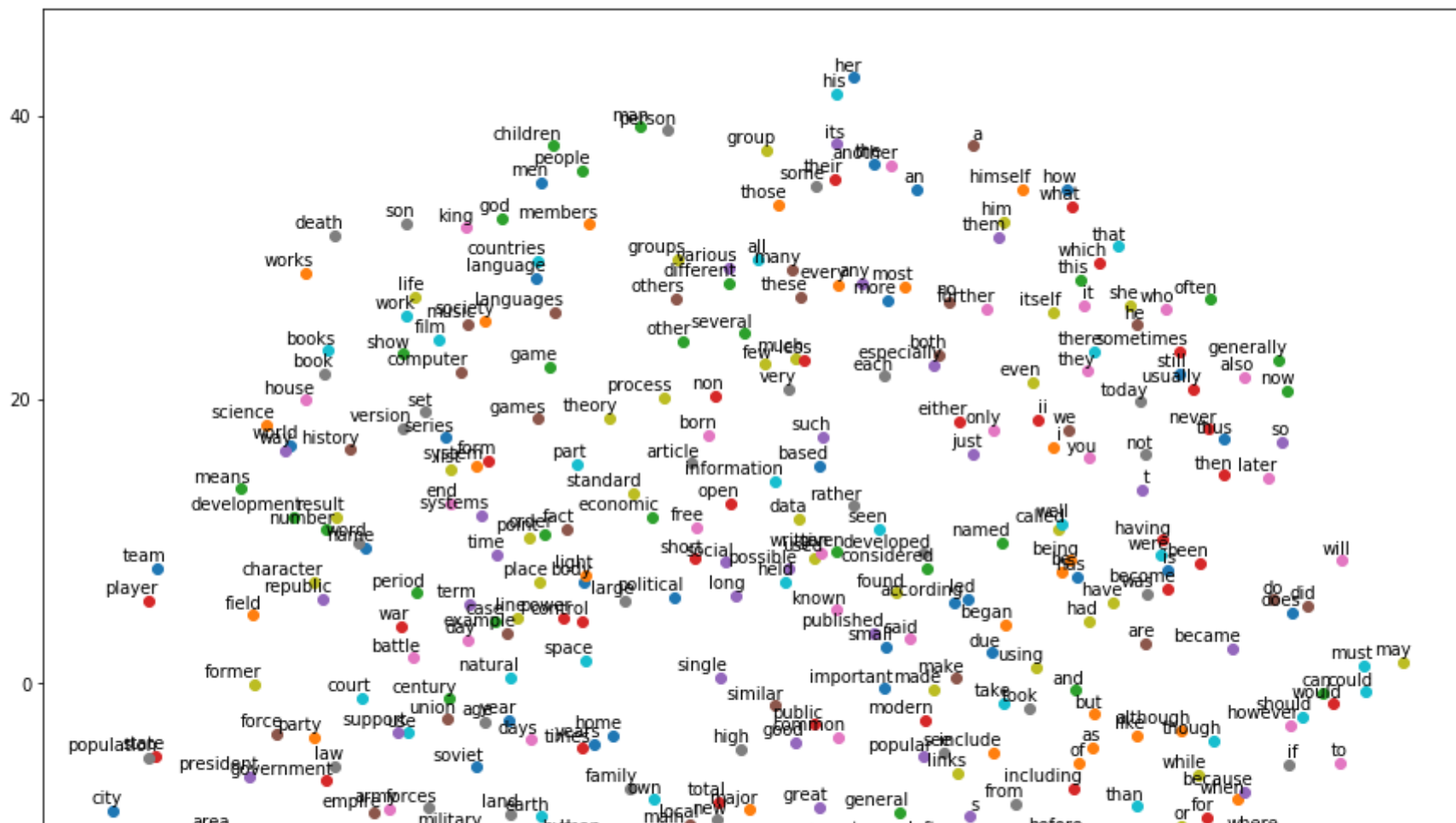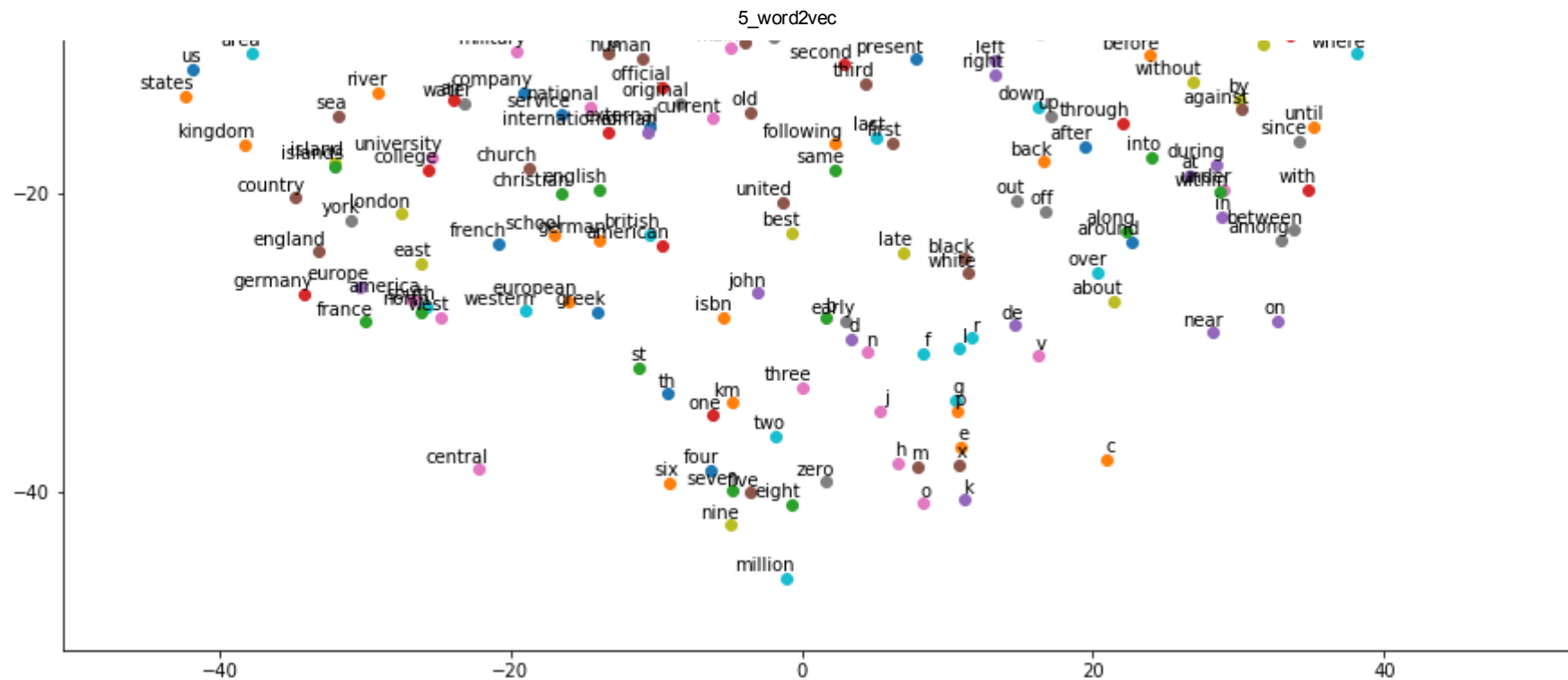
In [8]:
```
num_points = 400

tsne = TSNE(perplexity=30, n_components=2, init='pca', n_iter=5000)
two_d_embeddings = tsne.fit_transform(final_embeddings[1:num_points+1, :])
```

```
In [9]: def plot(embeddings, labels):
          assert embeddings.shape[0] >= len(labels), 'More labels than embeddings'
          pylab.figure(figsize=(15,15))  # in inches
          for i, label in enumerate(labels):
            x, y = embeddings[i,:]
            pylab.scatter(x, y)
            pylab.annotate(label, xy=(x, y), xytext=(5, 2), textcoords='offset points',
                           ha='right', va='bottom')
          pylab.show()

        words = [reverse_dictionary[i] for i in range(1, num_points+1)]
        plot(two_d_embeddings, words)
```

## Problem

An alternative to skip-gram is another Word2Vec model called CBOW (http://arxiv.org/abs/1301.3781) (Continuous Bag of Words). In the CBOW model, instead of predicting a context word from a word vector, you predict a word from the sum of all the word vectors in its context. Implement and evaluate a CBOW model trained on the text8 dataset.

```
In [11]: data_index = 0

         def generate_batch(batch_size, bag_window):
           global data_index
           span = 2 * bag_window + 1 # [ bag_window target bag_window ]
           batch = np.ndarray(shape=(batch_size, span - 1), dtype=np.int32)
           labels = np.ndarray(shape=(batch_size, 1), dtype=np.int32)
           buffer = collections.deque(maxlen=span)
           for _ in range(span):
             buffer.append(data[data_index])
             data_index = (data_index + 1) % len(data)
           for i in range(batch_size):
             # just for testing
             buffer_list = list(buffer)
             labels[i, 0] = buffer_list.pop(bag_window)
             batch[i] = buffer_list
             # iterate to the next buffer
             buffer.append(data[data_index])
             data_index = (data_index + 1) % len(data)
           return batch, labels

         print('data:', [reverse_dictionary[di] for di in data[:16]])

         for bag_window in [1, 2]:
           data_index = 0
           batch, labels = generate_batch(batch_size=4, bag_window=bag_window)
           print('\nwith bag_window = %d:' % (bag_window))
           print('    batch:', [[reverse_dictionary[w] for w in bi] for bi in batch])
           print('    labels:', [reverse_dictionary[li] for li in labels.reshape(4)])
```

```
data: ['anarchism', 'originated', 'as', 'a', 'term', 'of', 'abuse', 'first', 'used', 'against', 'early', 'working', 'cl
ass', 'radicals', 'including', 'the']

with bag_window = 1:
    batch: [['anarchism', 'as'], ['originated', 'a'], ['as', 'term'], ['a', 'of']]
    labels: ['originated', 'as', 'a', 'term']

with bag_window = 2:
    batch: [['anarchism', 'originated', 'a', 'term'], ['originated', 'as', 'term', 'of'], ['as', 'a', 'of', 'abuse'],
  ['a', 'term', 'abuse', 'first']]
    labels: ['as', 'a', 'term', 'of']
```

```
In [15]: batch_size = 128
         embedding_size = 128 # Dimension of the embedding vector.
         # skip_window = 1 # How many words to consider left and right.
         # num_skips = 2 # How many times to reuse an input to generate a label.
         bag_window = 2
         # We pick a random validation set to sample nearest neighbors. here we limit the
         # validation samples to the words that have a low numeric ID, which by
         # construction are also the most frequent.
         valid_size = 16 # Random set of words to evaluate similarity on.
         valid_window = 100 # Only pick dev samples in the head of the distribution.
         valid_examples = np.array(random.sample(range(valid_window), valid_size))
         num_sampled = 64 # Number of negative examples to sample.

         graph = tf.Graph()

         with graph.as_default(), tf.device('/cpu:0'):

           # Input data.
           train_dataset = tf.placeholder(tf.int32, shape=[batch_size, bag_window * 2])
           train_labels = tf.placeholder(tf.int32, shape=[batch_size, 1])
           valid_dataset = tf.constant(valid_examples, dtype=tf.int32)

           # Variables.
           embeddings = tf.Variable(
             tf.random_uniform([vocabulary_size, embedding_size], -1.0, 1.0))
           softmax_weights = tf.Variable(
             tf.truncated_normal([vocabulary_size, embedding_size],
                                 stddev=1.0 / math.sqrt(embedding_size)))
           softmax_biases = tf.Variable(tf.zeros([vocabulary_size]))

           # Model.
           # Look up embeddings for inputs.
           embed = tf.nn.embedding_lookup(embeddings, train_dataset)
           # Compute the softmax loss, using a sample of the negative labels each time.
           loss = tf.reduce_mean(
             tf.nn.sampled_softmax_loss(weights=softmax_weights, biases=softmax_biases, inputs=tf.reduce_sum(embed, 1),
                                        labels=train_labels, num_sampled=num_sampled, num_classes=vocabulary_size))

           # Optimizer.
           # Note: The optimizer will optimize the softmax_weights AND the embeddings.
           # This is because the embeddings are defined as a variable quantity and the
```

```
# optimizer's `minimize` method will by default modify all variable quantities
# that contribute to the tensor it is passed.
# See docs on `tf.train.Optimizer.minimize()` for more details.
optimizer = tf.train.AdagradOptimizer(1.0).minimize(loss)

# Compute the similarity between minibatch examples and all embeddings.
# We use the cosine distance:
norm = tf.sqrt(tf.reduce_sum(tf.square(embeddings), 1, keep_dims=True))
normalized_embeddings = embeddings / norm
valid_embeddings = tf.nn.embedding_lookup(
  normalized_embeddings, valid_dataset)
similarity = tf.matmul(valid_embeddings, tf.transpose(normalized_embeddings))
```

In [16]:
```python
num_steps = 100001

with tf.Session(graph=graph) as session:
  tf.global_variables_initializer().run()
  print('Initialized')
  average_loss = 0
  for step in range(num_steps):
    batch_data, batch_labels = generate_batch(batch_size, bag_window)
    feed_dict = {train_dataset : batch_data, train_labels : batch_labels}
    _, l = session.run([optimizer, loss], feed_dict=feed_dict)
    average_loss += l
    if step % 2000 == 0:
      if step > 0:
        average_loss = average_loss / 2000
      # The average loss is an estimate of the loss over the last 2000 batches.
      print('Average loss at step %d: %f' % (step, average_loss))
      average_loss = 0
    # note that this is expensive (~20% slowdown if computed every 500 steps)
    if step % 10000 == 0:
      sim = similarity.eval()
      for i in range(valid_size):
        valid_word = reverse_dictionary[valid_examples[i]]
        top_k = 8 # number of nearest neighbors
        nearest = (-sim[i, :]).argsort()[1:top_k+1]
        log = 'Nearest to %s:' % valid_word
        for k in range(top_k):
          close_word = reverse_dictionary[nearest[k]]
          log = '%s %s,' % (log, close_word)
        print(log)
  final_embeddings = normalized_embeddings.eval()
```

```
Nearest to time: distance, period, least, point, thing, looking, year, night,
Nearest to if: when, though, where, subcategories, therefore, since, currently, whatever,
Nearest to when: if, where, until, although, while, though, since, because,
Nearest to on: upon, through, in, amanda, against, across, otomo, fogo,
Nearest to while: though, and, although, when, amongst, or, where, however,
Nearest to state: government, moran, lecter, observant, cumberland, session, copy, annexes,
Nearest to of: following, almohades, despite, including, amputee, discernable, regarding, otomo,
Nearest to with: between, tuba, via, among, contended, terminus, containing, anoint,
Nearest to who: benn, always, young, sorties, raps, whom, which, attenborough,
Nearest to not: never, indeed, still, nothing, legally, almost, also, cyclopedia,
Nearest to by: when, through, without, pelvic, in, sopranos, crocodile, actually,
```
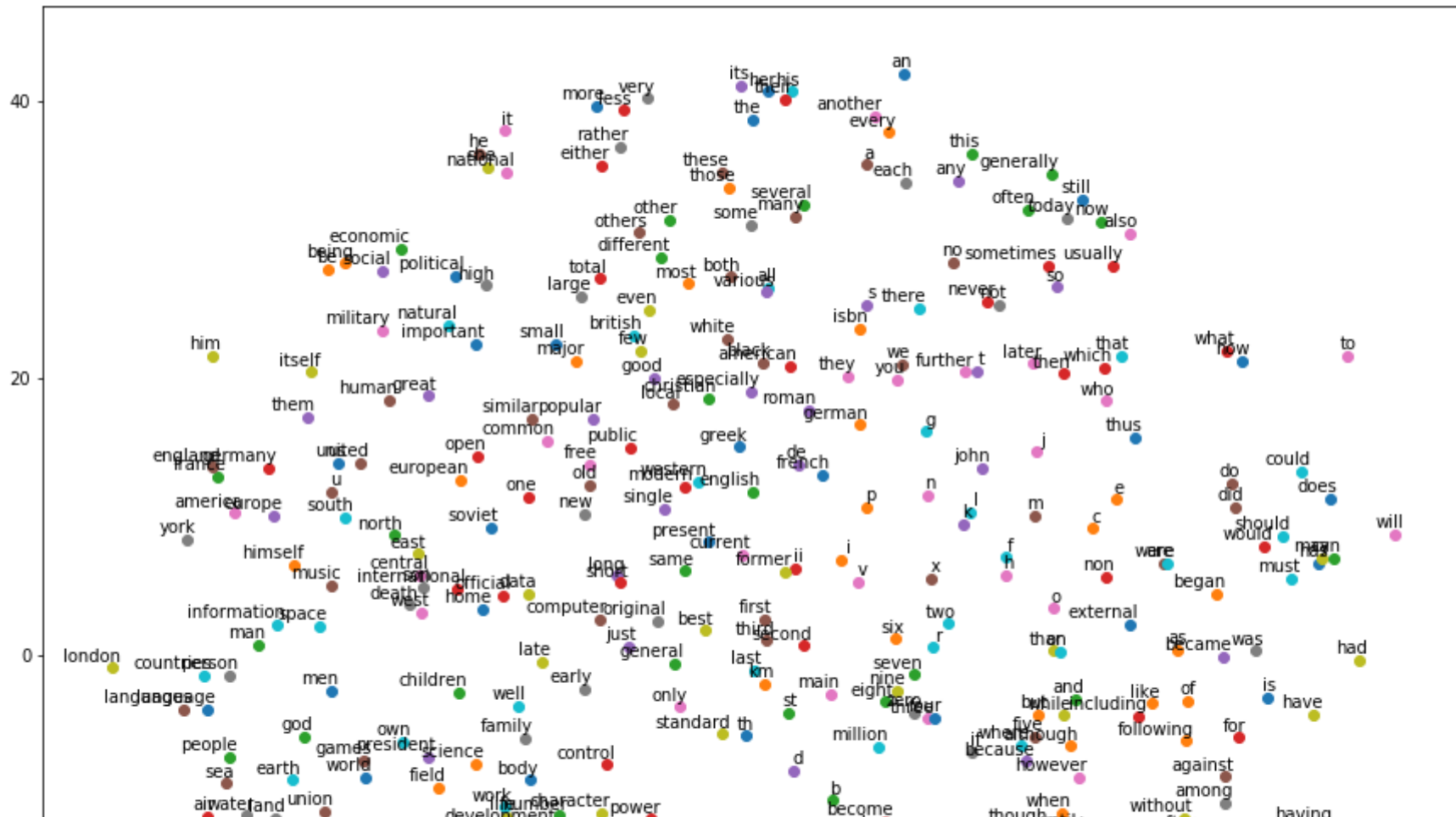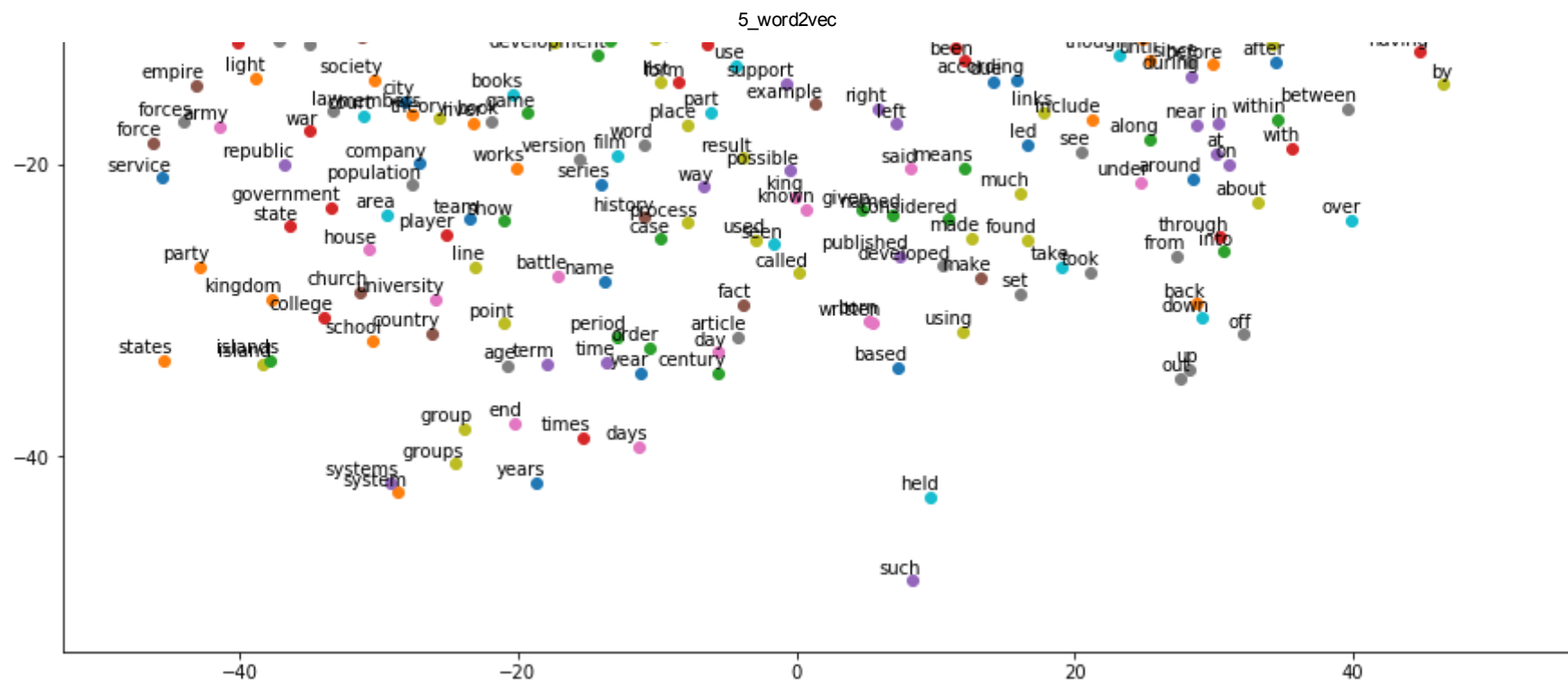
```
Nearest to other: others, various, textiles, different, both, roms, hae, fewer,
Nearest to states: kingdom, nations, emirates, felt, coordinating, countries, yorkist, declared,
Average loss at step 92000: 2.894091
Average loss at step 94000: 2.882162
Average loss at step 96000: 2.723037
Average loss at step 98000: 2.457204
Average loss at step 100000: 2.715885
Nearest to can: could, must, should, cannot, might, may, will, would,
Nearest to its: their, his, her, the, our, your, patrilineal, abductions,
```

In [17]:
```python
num_points = 400

tsne = TSNE(perplexity=30, n_components=2, init='pca', n_iter=5000)
two_d_embeddings = tsne.fit_transform(final_embeddings[1:num_points+1, :])
```

```
In [18]: def plot(embeddings, labels):
           assert embeddings.shape[0] >= len(labels), 'More labels than embeddings'
           pylab.figure(figsize=(15,15))  # in inches
           for i, label in enumerate(labels):
             x, y = embeddings[i,:]
             pylab.scatter(x, y)
             pylab.annotate(label, xy=(x, y), xytext=(5, 2), textcoords='offset points',
                            ha='right', va='bottom')
           pylab.show()

         words = [reverse_dictionary[i] for i in range(1, num_points+1)]
         plot(two_d_embeddings, words)
```

In [ ]: