Assignment 5

The goal of this assignment is to train a Word2Vec skip-gram model over Text8 data. In [1]: # 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]: 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]: 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]
      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</pre>
  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 ]
```

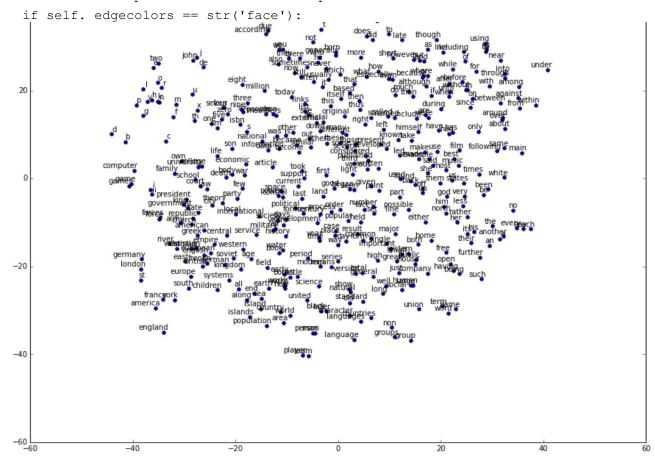
```
for j in range(num_skips):
      while target in targets_to_avoid:
        target = random.randint(0, span - 1)
      targets_to_avoid.append(target)
      batch[i * num_skips + j] = buffer[skip_window]
      labels[i * num_skips + j, 0] = buffer[target]
    buffer.append(data[data_index])
    data_index = (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: ['as', 'anarchism', 'originated', 'a', 'as', 'term', 'of', 'a']
with num_skips = 4 and skip_window = 2:
    batch: ['as', 'as', 'as', 'a', 'a', 'a', 'a']
    labels: ['term', 'a', 'originated', 'anarchism', 'term', 'of',
'originated', '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]))
  # Mode1.
  # 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(softmax_weights, softmax_biases, embed,
                               train_labels, num_sampled, 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]:
num\_steps = 100001
with tf.Session(graph=graph) as session:
  tf.initialize_all_variables().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 += 1
    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()
Initialized
Average loss at step 0: 8.161867
Nearest to all: ajaccio, legged, creationism, vagaries, guitarists, penetration,
ibero, mactutor,
Nearest to history: diaeresis, nicolai, canmore, mudslides, harmonize, excerpt,
thujone, fried,
Nearest to from: eclac, texan, subsisting, kramer, arminians, traction,
displayed, dissatisfied,
Nearest to system: arthur, gmt, midair, zug, thicker, pathological, fabricated,
Nearest to or: omnium, condemns, lama, keen, capitulated, participate, removes,
aglaulus,
Nearest to often: sometimes, usually, commonly, generally, frequently, now,
Nearest to four: six, seven, eight, five, three, nine, two, zero,
Nearest to new: different, painterly, huge, masoretic, mohiam, scuba, outback,
old,
Nearest to there: they, it, he, we, still, often, usually, she,
                                                                          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)
```

/home/maxkhk/anaconda/lib/python2.7/site-packages/matplotlib/collections.py:590: FutureWarning: elementwise comparison failed; returning scalar instead, but in the future will perform elementwise comparison



Problem

An alternative to skip-gram is another Word2Vec model called <u>CBOW</u> (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 [9]:
```

```
# 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
from itertools import compress
```

```
#data preparation
```

```
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)
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))
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.
Found and verified text8.zip
Data size 17005207
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]
                                                                         In [20]:
data_index = 0
#generally what we do here is build data as a pack of words to left and right
of the target word
#and label is the target word
def generate_batch_cbow(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_cbow(batch_size=4, bag_window=bag_window)
  print('\nwith bag_window = %d:' % (bag_window))
             batch:', [[reverse_dictionary[w] for w in bi] for bi in batch])
  print('
  print('
             labels:', [reverse_dictionary[li] for li in labels.reshape(4)])
data: ['anarchism', 'originated', 'as', 'a', 'term', 'of', 'abuse', 'first',
'used', 'against', 'early', 'working', 'class', '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 [21]:
batch size = 128
embedding_size = 128 # Dimension of the embedding vector.
bag_window = 2 # How many words to consider left and right.
# We pick a random validation set to sample nearest neighbors. here we limit
# 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]))
  # Mode1.
  # Look up embeddings for inputs.
  embeds = tf.nn.embedding_lookup(embeddings, train_dataset)
  # Compute the softmax loss, using a sample of the negative labels each time.
  # Notice per description of CBOW we have to sum all the embedings!
  loss = tf.reduce_mean(
    tf.nn.sampled_softmax_loss(softmax_weights, softmax_biases,
tf.reduce_sum(embeds,1),
                               train_labels, num_sampled, vocabulary_size))
  # Optimizer.
  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))
num\_steps = 100001
with tf.Session(graph=graph) as session:
  tf.initialize_all_variables().run()
  print('Initialized')
  average_loss = 0
  for step in range(num_steps):
    batch_data, batch_labels = generate_batch_cbow(
      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 += 1
    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()
Initialized
Average loss at step 0: 8.095695
Nearest to people: ecuador, visibly, contraceptives, bitstream, howell, chagrin,
launcher, circumvented,
Nearest to this: shimura, niki, squealer, tachycardia, sketchy, solubility,
sheer, australis,
Nearest to in: hotly, raions, expended, papen, lachlan, tis, calhoun, stipend,
Nearest to of: mostar, rudolph, divert, berthold, babel, sewage, agi, antony,
Nearest to many: some, several, various, numerous, both, all, hundreds, most,
Nearest to their: its, his, her, your, our, my, the, whose,
Nearest to and: or, but, while, including, however, although, stumble, etc,
Nearest to the: its, his, their, a, any, every, our, each,
```

/home/maxkhk/anaconda/lib/python2.7/site-packages/matplotlib/collections.py:590: FutureWarning: elementwise comparison failed; returning scalar instead, but in the future will perform elementwise comparison

if self._edgecolors == str('face'):

