

# Skip-gram word2vec

In this notebook, I'll lead you through using TensorFlow to implement the word2vec algorithm using the skip-gram architecture. By implementing this, you'll learn about embedding words for use in natural language processing. This will come in handy when dealing with things like translations.

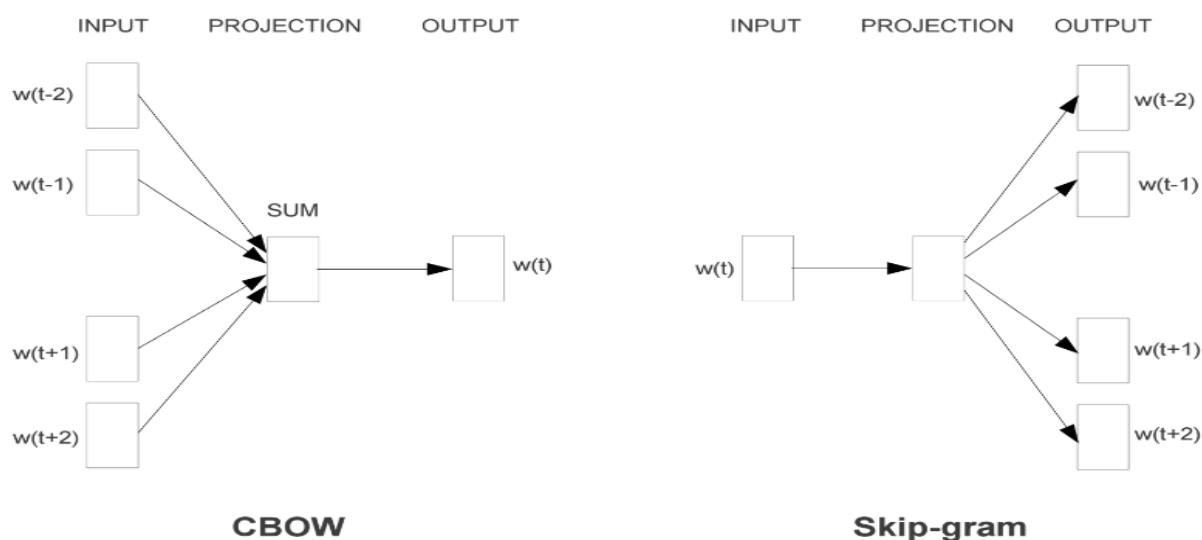
## Readings

Here are the resources I used to build this notebook. I suggest reading these either beforehand or while you're working on this material.

- A really good [conceptual overview](#) of word2vec from Chris McCormick
- [First word2vec paper](#) from Mikolov et al.
- [NIPS paper](#) with improvements for word2vec also from Mikolov et al.
- An [implementation of word2vec](#) from Thushan Ganegedara
- TensorFlow [word2vec tutorial](#)

## Word embeddings

When you're dealing with language and words, you end up with tens of thousands of classes to predict, one for each word. Trying to one-hot encode these words is massively inefficient, you'll have one element set to 1 and the other 50,000 set to 0. The word2vec algorithm finds much more efficient representations by finding vectors that represent the words. These vectors also contain semantic information about the words. Words that show up in similar contexts, such as "black", "white", and "red" will have vectors near each other. There are two architectures for implementing word2vec, CBOW (Continuous Bag-Of-Words) and Skip-gram.



In this implementation, we'll be using the skip-gram architecture because it performs better than CBOW. Here, we pass in a word and try to predict the words surrounding it in the text. In this way, we can train the network to learn representations for words that show up in similar contexts. First up, importing packages.

In [1]:

```
import time
```

```
import numpy as np
import tensorflow as tf
```

```
import utils
```

Load the [text8 dataset](#), a file of cleaned up Wikipedia articles from Matt Mahoney. The next cell will download the data set to the `data` folder. Then you can extract it and delete the archive file to save storage space.

In [2]:

```

from urllib.request import urlretrieve
from os.path import isfile, isdir
from tqdm import tqdm
import zipfile

dataset_folder_path = 'data'
dataset_filename = 'text8.zip'
dataset_name = 'Text8 Dataset'

class DLProgress(tqdm):
    last_block = 0

    def hook(self, block_num=1, block_size=1, total_size=None):
        self.total = total_size
        self.update((block_num - self.last_block) * block_size)
        self.last_block = block_num

if not isfile(dataset_filename):
    with DLProgress(unit='B', unit_scale=True, miniters=1, desc=dataset_name)
    as pbar:
        urlretrieve(
            'http://mattmahoney.net/dc/text8.zip',
            dataset_filename,
            pbar.hook)

if not isdir(dataset_folder_path):
    with zipfile.ZipFile(dataset_filename) as zip_ref:
        zip_ref.extractall(dataset_folder_path)

with open('data/text8') as f:
    text = f.read()

```

---

Text8 Dataset: 31.4MB [00:16, 1.88MB/s]

## Preprocessing

Here I'm fixing up the text to make training easier. This comes from the `utils` module I wrote. The `preprocess` function converts any punctuation into tokens, so a period is changed to `<PERIOD>`. In this data set, there aren't any periods, but it will help in other NLP problems. I'm also removing all words that show up five or fewer times in the dataset. This will greatly reduce issues due to noise in the data and improve the quality of the vector representations. If you want to write your own functions for this stuff, go for it.

In [57]:

```

words = utils.preprocess(text)
print(words[:30])
['anarchism', 'originated', 'as', 'a', 'term', 'of', 'abuse', 'first', 'used',
'against', 'early', 'working', 'class', 'radicals', 'including', 'the',
'diggers', 'of', 'the', 'english', 'revolution', 'and', 'the', 'sans',
'culottes', 'of', 'the', 'french', 'revolution', 'whilst']

```

In [58]:

```

print("Total words: {}".format(len(words)))
print("Unique words: {}".format(len(set(words))))
Total words: 16680599
Unique words: 63641

```

And here I'm creating dictionaries to covert words to integers and backwards, integers to words. The integers are assigned in descending frequency order, so the most frequent word ("the") is given the integer 0 and the next most frequent is 1 and so on. The words are converted to integers and stored in the list `int_words`.

In [59]:

```
vocab_to_int, int_to_vocab = utils.create_lookup_tables(words)
int_words = [vocab_to_int[word] for word in words]
```

## Subsampling

Words that show up often such as "the", "of", and "for" don't provide much context to the nearby words. If we discard some of them, we can remove some of the noise from our data and in return get faster training and better representations. This process is called subsampling by Mikolov. For each word  $w_i$  in the training set, we'll discard it with probability given by

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

where  $t$  is a threshold parameter and  $f(w_i)$  is the frequency of word  $w_i$  in the total dataset.

I'm going to leave this up to you as an exercise. Check out my solution to see how I did it.

**Exercise:** Implement subsampling for the words in `int_words`. That is, go through `int_words` and discard each word given the probability  $P(w_i)$  shown above. Note that  $P(w_i)$  is that probability that a word is discarded. Assign the subsampled data to `train_words`.

In [60]:

```
from collections import Counter
import random

threshold = 1e-5
word_counts = Counter(int_words)
total_count = len(int_words)
freqs = {word: count/total_count for word, count in word_counts.items()}
p_drop = {word: 1 - np.sqrt(threshold/freqs[word]) for word in word_counts}
train_words = [word for word in int_words if p_drop[word] < random.random()]
```

## Making batches

Now that our data is in good shape, we need to get it into the proper form to pass it into our network. With the skip-gram architecture, for each word in the text, we want to grab all the words in a window around that word, with size  $CC$ .

From [Mikolov et al.](#):

"Since the more distant words are usually less related to the current word than those close to it, we give less weight to the distant words by sampling less from those words in our training examples... If we choose  $C=5$ , for each training word we will select randomly a number  $RR$  in range  $<1;C>$ , and then use  $RR$  words from history and  $RR$  words from the future of the current word as correct labels."

**Exercise:** Implement a function `get_target` that receives a list of words, an index, and a window size, then returns a list of words in the window around the index. Make sure to use the algorithm described above, where you chose a random number of words to from the window.

In [61]:

```
def get_target(words, idx, window_size=5):
    ''' Get a list of words in a window around an index. '''

    R = np.random.randint(1, window_size+1)
    start = idx - R if (idx - R) > 0 else 0
    stop = idx + R
    target_words = set(words[start:idx] + words[idx+1:stop+1])

    return list(target_words)
```

Here's a function that returns batches for our network. The idea is that it grabs `batch_size` words from a words list. Then for each of those words, it gets the target words in the window. I haven't found a way to pass in a random number of target words and get it to work with the architecture, so I make one row per input-target pair. This is a generator function by the way, helps save memory.

In [62]:

```
def get_batches(words, batch_size, window_size=5):
    ''' Create a generator of word batches as a tuple (inputs, targets) '''

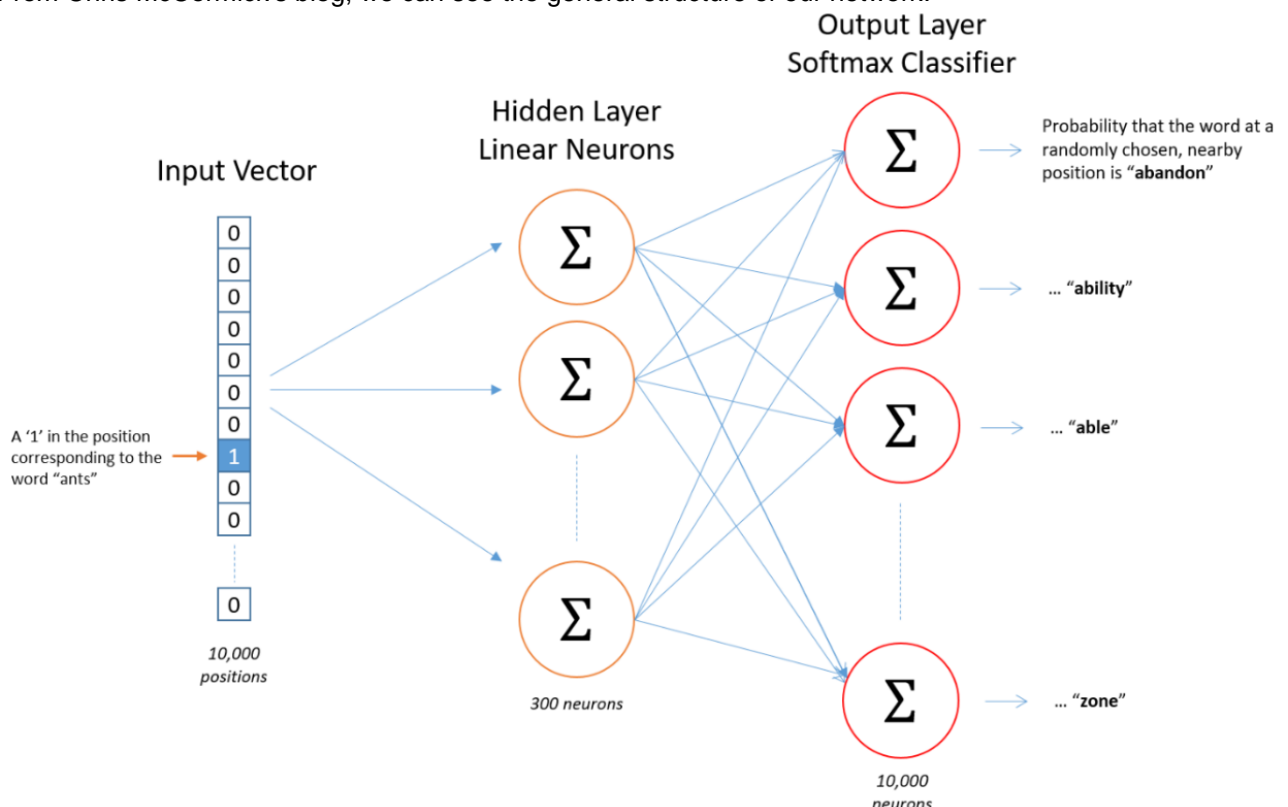
    n_batches = len(words)//batch_size

    # only full batches
    words = words[:n_batches*batch_size]

    for idx in range(0, len(words), batch_size):
        x, y = [], []
        batch = words[idx:idx+batch_size]
        for ii in range(len(batch)):
            batch_x = batch[ii]
            batch_y = get_target(batch, ii, window_size)
            y.extend(batch_y)
            x.extend([batch_x]*len(batch_y))
        yield x, y
```

## Building the graph

From Chris McCormick's blog, we can see the general structure of our network.



The input words are passed in as one-hot encoded vectors. This will go into a hidden layer of linear units, then into a softmax layer. We'll use the softmax layer to make a prediction like normal.

The idea here is to train the hidden layer weight matrix to find efficient representations for our words. This weight matrix is usually called the embedding matrix or embedding look-up table. We can discard the softmax layer because we don't really care about making predictions with this network. We just want the embedding matrix so we can use it in other networks we build from the dataset.

I'm going to have you build the graph in stages now. First off, creating the `inputs` and `labels` placeholders like normal.

**Exercise:** Assign `inputs` and `labels` using `tf.placeholder`. We're going to be passing in integers, so set the data types to `tf.int32`. The batches we're passing in will have varying sizes, so set the batch sizes to `[None]`. To make things work later, you'll need to set the second dimension of `labels` to `None` or `1`.

In [ ]:

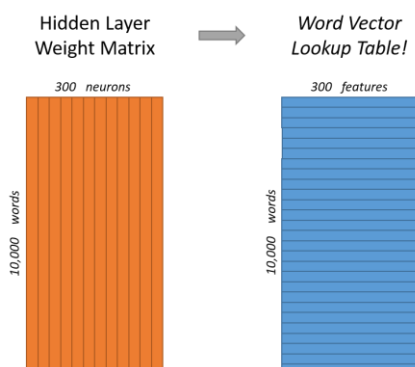
```
train_graph = tf.Graph()
with train_graph.as_default():
    inputs = tf.placeholder(tf.int32, [None], name='inputs')
    labels = tf.placeholder(tf.int32, [None, None], name='labels')
```

## Embedding

The embedding matrix has a size of the number of words by the number of neurons in the hidden layer. So, if you have 10,000 words and 300 hidden units, the matrix will have size 10,000×300. Remember that we're using one-hot encoded vectors for our inputs. When you do the matrix multiplication of the one-hot vector with the embedding matrix, you end up selecting only one row out of the entire matrix:

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

You don't actually need to do the matrix multiplication, you just need to select the row in the embedding matrix that corresponds to the input word. Then, the embedding matrix becomes a lookup table, you're looking up a vector the size of the hidden layer that represents the input word.



**Exercise:** Tensorflow provides a convenient function

`tf.nn.embedding_lookup` that does this lookup for us. You pass in the embedding matrix and a tensor of integers, then it returns rows in the matrix corresponding to those integers. Below, set the number of embedding features you'll use (200 is a good start), create the embedding matrix variable, and use `tf.nn.embedding_lookup` to get the embedding tensors. For the embedding matrix, I suggest you initialize it with a uniform random numbers between -1 and 1 using [tf.random\\_uniform](#).

In [ ]:

```
n_vocab = len(int_to_vocab)
n_embedding = 200 # Number of embedding features
with train_graph.as_default():
    embedding = tf.Variable(tf.random_uniform((n_vocab, n_embedding), -1, 1))
    embed = tf.nn.embedding_lookup(embedding, inputs)
```

## Negative sampling

For every example we give the network, we train it using the output from the softmax layer. That means for each input, we're making very small changes to millions of weights even though we only have one true example. This makes training the network very inefficient. We can approximate the loss from the softmax layer by only updating a small subset of all the weights at once. We'll update the weights for the correct label, but only a small number of incorrect labels. This is called ["negative sampling"](#). Tensorflow has a convenient function to do this, `tf.nn.sampled_softmax_loss`.

**Exercise:** Below, create weights and biases for the softmax layer. Then, use `tf.nn.sampled_softmax_loss` to calculate the loss. Be sure to read the documentation to figure out how it works.

In [66]:

```

# Number of negative labels to sample
n_sampled = 100
with train_graph.as_default():
    softmax_w = tf.Variable(tf.truncated_normal((n_vocab, n_embedding),
stddev=0.1))
    softmax_b = tf.Variable(tf.zeros(n_vocab))

    # Calculate the loss using negative sampling
    loss = tf.nn.sampled_softmax_loss(softmax_w, softmax_b,
                                      labels, embed,
                                      n_sampled, n_vocab)

    cost = tf.reduce_mean(loss)
    optimizer = tf.train.AdamOptimizer().minimize(cost)

```

## Validation

This code is from Thushan Ganegedara's implementation. Here we're going to choose a few common words and few uncommon words. Then, we'll print out the closest words to them. It's a nice way to check that our embedding table is grouping together words with similar semantic meanings.

In [ ]:

```

with train_graph.as_default():
    ## From Thushan Ganegedara's implementation
    valid_size = 16 # Random set of words to evaluate similarity on.
    valid_window = 100
    # pick 8 samples from (0,100) and (1000,1100) each ranges. lower id implies
more frequent
    valid_examples = np.array(random.sample(range(valid_window), valid_size//2))
    valid_examples = np.append(valid_examples,
                              random.sample(range(1000,1000+valid_window), valid_size//2))

    valid_dataset = tf.constant(valid_examples, dtype=tf.int32)

    # We use the cosine distance:
    norm = tf.sqrt(tf.reduce_sum(tf.square(embedding), 1, keep_dims=True))
    normalized_embedding = embedding / norm
    valid_embedding = tf.nn.embedding_lookup(normalized_embedding, valid_dataset)
    similarity = tf.matmul(valid_embedding, tf.transpose(normalized_embedding))

```

In [18]:

```

# If the checkpoints directory doesn't exist:
!mkdir checkpoints

```

In [67]:

```

epochs = 10
batch_size = 1000
window_size = 10

with train_graph.as_default():
    saver = tf.train.Saver()

with tf.Session(graph=train_graph) as sess:
    iteration = 1

```

```

loss = 0
sess.run(tf.global_variables_initializer())

for e in range(1, epochs+1):
    batches = get_batches(train_words, batch_size, window_size)
    start = time.time()
    for x, y in batches:

        feed = {inputs: x,
                labels: np.array(y)[: , None]}
        train_loss, _ = sess.run([cost, optimizer], feed_dict=feed)

        loss += train_loss

    if iteration % 100 == 0:
        end = time.time()
        print("Epoch {}/{}".format(e, epochs),
              "Iteration: {}".format(iteration),
              "Avg. Training loss: {:.4f}".format(loss/100),
              "{:.4f} sec/batch".format((end-start)/100))
        loss = 0
        start = time.time()

    if iteration % 1000 == 0:
        # note that this is expensive (~20% slowdown if computed every 500 steps)
        sim = similarity.eval()

        for i in range(valid_size):
            valid_word = int_to_vocab[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 = int_to_vocab[nearest[k]]
                log = '%s %s,' % (log, close_word)
            print(log)

        iteration += 1
    save_path = saver.save(sess, "checkpoints/text8.ckpt")
    embed_mat = sess.run(normalized_embedding)

```

---

Epoch 1/10 Iteration: 100 Avg. Training loss: 5.6559 0.1018 sec/batch

Epoch 1/10 Iteration: 200 Avg. Training loss: 5.6093 0.1028 sec/batch

.....

Nearest to channel: creditors, channels, curler, hearsay, mbit, wb, carnivores,  
bandwidth,

Nearest to report: reports, credibility, annotated, commission, zangger, santer,  
focusing, lists,

Epoch 10/10 Iteration: 46100 Avg. Training loss: 3.8255 0.1184 sec/batch

Epoch 10/10 Iteration: 46200 Avg. Training loss: 3.8518 0.1119 sec/batch



Restore the trained network if you need to:

In [20]:

```
with train_graph.as_default():
    saver = tf.train.Saver()

with tf.Session(graph=train_graph) as sess:
    saver.restore(sess, tf.train.latest_checkpoint('checkpoints'))
    embed_mat = sess.run(embedding)
```

## Visualizing the word vectors

Below we'll use T-SNE to visualize how our high-dimensional word vectors cluster together. T-SNE is used to project these vectors into two dimensions while preserving local structure. Check out [this post from Christopher Olah](#) to learn more about T-SNE and other ways to visualize high-dimensional data.

In [115]:

```
%matplotlib inline
%config InlineBackend.figure_format = 'retina'
```

```
import matplotlib.pyplot as plt
from sklearn.manifold import TSNE
```

In [138]:

```
viz_words = 500
tsne = TSNE()
embed_tsne = tsne.fit_transform(embed_mat[:viz_words, :])
```

In [139]:

```
fig, ax = plt.subplots(figsize=(14, 14))
for idx in range(viz_words):
    plt.scatter(*embed_tsne[idx, :], color='steelblue')
    plt.annotate(int_to_vocab[idx], (embed_tsne[idx, 0], embed_tsne[idx, 1]),
alpha=0.7)
```

