

⊖ Everything running smoothly!

Applying Recurrent Neural Networks/LSTM for Language Modeling

Hello and welcome to this part. In this notebook, we will go over the topic of Language Modelling, and create a Recurrent Neural Network model based on the Long Short-Term Memory unit to train and benchmark on the Penn Treebank dataset. By the end of this notebook, you should be able to understand how TensorFlow builds and executes a RNN model for Language Modelling.

The Objective

By now, you should have an understanding of how Recurrent Networks work -- a specialized model to process sequential data by keeping track of the "state" or context. In this notebook, we go over a TensorFlow code snippet for creating a model focused on **Language Modelling** -- a very relevant task that is the cornerstone of many different linguistic problems such as **Speech Recognition, Machine Translation and Image Captioning**. For this, we will be using the Penn Treebank dataset, which is an often-used dataset for benchmarking Language Modelling models.

Table of Contents

1. [What exactly is Language Modelling?](#)
2. [The Penn Treebank dataset](#)
3. [Word Embedding](#)
4. [Building the LSTM model for Language Modeling](#)
5. [LSTM](#)

</div>

What exactly is Language Modelling?

Language Modelling, to put it simply, **is the task of assigning probabilities to sequences of words**. This means that, given a context of one or a sequence of words in the language the model was trained on, the model should provide the next most probable words or sequence of words that follows from the given sequence of words the sentence. Language Modelling is one of the most important tasks in Natural Language Processing.

In this example, one can see the predictions for the next word of a sentence, given the context "This is an". As you can see, this boils down to a sequential data analysis task -- you are given a word or a sequence of words (the input data), and, given the context (the state), you need to find out what is the next word (the prediction). This kind of analysis is very important for language-related tasks such as **Speech Recognition, Machine Translation, Image Captioning, Text Correction** and many other very relevant problems.

As the above image shows, Recurrent Network models fit this problem like a glove. Alongside LSTM and its capacity to maintain the model's state for over one thousand time steps, we have all the tools we need to undertake this problem. The goal for this notebook is to create a model that can reach **low levels of perplexity** on our desired dataset.

For Language Modelling problems, **perplexity** is the way to gauge efficiency. Perplexity is simply a measure of how well a probabilistic model is able to predict its sample. A higher-level way to explain this would be saying that **low perplexity means a higher degree of trust in the predictions the model makes**. Therefore, the lower perplexity is, the better.

The Penn Treebank dataset

Historically, datasets big enough for Natural Language Processing are hard to come by. This is in part due to the necessity of the sentences to be broken down and tagged with a certain degree of correctness -- or else the models trained on it won't be able to be correct at all. This means that we need a **large amount of data, annotated by or at least corrected by humans**. This is, of course, not an easy task at all.

The Penn Treebank, or PTB for short, is a dataset maintained by the University of Pennsylvania. It is *huge* -- there are over **four million and eight hundred thousand** annotated words in it, all corrected by humans. It is composed of many different sources, from abstracts of Department of Energy papers to texts from the Library of America. Since it is verifiably correct and of such a huge size, the Penn Treebank is commonly used as a benchmark dataset for Language Modelling.

The dataset is divided in different kinds of annotations, such as Piece-of-Speech, Syntactic and Semantic skeletons. For this example, we will simply use a sample of clean, non-annotated words (with the exception of one tag --<unk> , which is used for rare words such as uncommon proper nouns) for our model. This means that we just want to predict what the next words would be, not what they mean in context or their classes on a given sentence.

Word Embeddings

For better processing, in this example, we will make use of [word embeddings](#), which is **a way of representing sentence structures or words as n-dimensional vectors (where n is a reasonably high number, such as 200 or 500) of real numbers**. Basically, we will assign each word a randomly-initialized vector, and input those into the network to be processed. After a number of iterations, these vectors are expected to assume values that help the network to correctly predict what it needs to -- in our case, the probable next word in the sentence. This is shown to be a very effective task in Natural Language

Processing, and is a commonplace practice.

Word Embedding tends to group up similarly used words *reasonably* close together in the vectorial space. For example, if we use T-SNE (a dimensional reduction visualization algorithm) to flatten the dimensions of our vectors into a 2-dimensional space and plot these words in a 2-dimensional space, we might see something like this:

As you can see, words that are frequently used together, in place of each other, or in the same places as them tend to be grouped together -- being closer together the higher they are correlated. For example, "None" is pretty semantically close to "Zero", while a phrase that uses "Italy", you could probably also fit "Germany" in it, with little damage to the sentence structure. The vectorial "closeness" for similar words like this is a great indicator of a well-built model.

We need to import the necessary modules for our code. We need **numpy** and **tensorflow**, obviously. Additionally, we can import directly the **tensorflow.models.rnn** model, which includes the function for building RNNs, and **tensorflow.models.rnn.ptb.reader** which is the helper module for getting the input data from the dataset we just downloaded.

If you want to learn more take a look at <https://github.com/tensorflow/models/blob/master/tutorials/rnn/ptb/reader.py>

In [7]:

```
import time
import numpy as np
import tensorflow as tf
import reader
```

In [7]:

```
!mkdir data
!wget -q -O data/ptb.zip https://ibm.box.com/shared/static/
z2yvmhbskc45xd2a9a4kkn6hg4g4kj5r.zip
!unzip -o data/ptb.zip -d data
!cp data/ptb/reader.py .
```

mkdir: cannot create directory 'data': File exists

Archive: data/ptb.zip

inflating: data/ptb/reader.py

inflating: data/___MACOSX/ptb/._reader.py

```
inflating: data/__MACOSX/._ptb
```

Building the LSTM model for Language Modeling

Now that we know exactly what we are doing, we can start building our model using TensorFlow. The very first thing we need to do is download and extract the `simple-examples` dataset, which can be done by executing the code cell below.

In [8]:

```
!wget http://www.fit.vutbr.cz/~imikolov/rnnlm/simple-
examples.tgz
!tar xzf simple-examples.tgz -C data/
```

```
--2020-02-02 11:28:39-- http://www.fit.vutbr.cz/~imikolov/
rnnlm/simple-examples.tgz
Resolving www.fit.vutbr.cz (www.fit.vutbr.cz)...
147.229.9.23, 2001:67c:1220:809::93e5:917
Connecting to www.fit.vutbr.cz (www.fit.vutbr.cz)|
147.229.9.23|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 34869662 (33M) [application/x-gtar]
Saving to: 'simple-examples.tgz.1'
```

```
simple-examples.tgz 100%[=====>] 33.25M
4.10MB/s in 9.3s
```

```
2020-02-02 11:28:49 (3.56 MB/s) - 'simple-examples.tgz.1'
saved [34869662/34869662]
```

Additionally, for the sake of making it easy to play around with the model's hyperparameters, we can declare them beforehand. Feel free to change these -- you will see a difference in performance each time you change those!

In [8]:

```
#Initial weight scale
init_scale = 0.1
#Initial learning rate
learning_rate = 1.0
#Maximum permissible norm for the gradient (For gradient
clipping -- another measure against Exploding Gradients)
max_grad_norm = 5
#The number of layers in our model
num_layers = 2
```

```

#The total number of recurrence steps, also known as the
number of layers when our RNN is "unfolded"
num_steps = 20
#The number of processing units (neurons) in the hidden
layers
hidden_size_l1 = 256
hidden_size_l2 = 128
#The maximum number of epochs trained with the initial
learning rate
max_epoch_decay_lr = 4
#The total number of epochs in training
max_epoch = 15
#The probability for keeping data in the Dropout Layer
(This is an optimization, but is outside our scope for this
notebook!)
#At 1, we ignore the Dropout Layer wrapping.
keep_prob = 1
#The decay for the learning rate
decay = 0.5
#The size for each batch of data
batch_size = 60
#The size of our vocabulary
vocab_size = 10000
embedding_vector_size = 200
#Training flag to separate training from testing
is_training = 1
#Data directory for our dataset
data_dir = "data/simple-examples/data/"

```

Some clarifications for LSTM architecture based on the arguments:

Network structure:

- In this network, the number of LSTM cells are 2. To give the model more expressive power, we can add multiple layers of LSTMs to process the data. The output of the first layer will become the input of the second and so on.
- The recurrence steps is 20, that is, when our RNN is "Unfolded", the recurrence step is 20.
- the structure is like:
 - 200 input units -> [200x200] Weight -> 200 Hidden units (first layer) -> [200x200] Weight matrix -> 200 Hidden units (second layer) -> [200] weight Matrix -> 200 unit output
- Input layer:

- The network has 200 input units.
- Suppose each word is represented by an embedding vector of dimensionality $e=200$. The input layer of each cell will have 200 linear units. These $e=200$ linear units are connected to each of the $h=200$ LSTM units in the hidden layer (assuming there is only one hidden layer, though our case has 2 layers).
- The input shape is $[\text{batch_size}, \text{num_steps}]$, that is $[30 \times 20]$. It will turn into $[30 \times 20 \times 200]$ after embedding, and then $20 \times [30 \times 200]$

•

Hidden layer:

- Each LSTM has 200 hidden units which is equivalent to the dimensionality of the embedding words and output.

•

There is a lot to be done and a ton of information to process at the same time, so go over this code slowly. It may seem complex at first, but if you try to apply what you just learned about language modelling to the code you see, you should be able to understand it.

This code is adapted from the [PTBModel](#) example bundled with the TensorFlow source code.

Train data

The story starts from data:

- Train data is a list of words, of size 929589, represented by numbers, e.g. [9971, 9972, 9974, 9975,...]
- We read data as mini-batch of size $b=30$. Assume the size of each sentence is 20 words ($\text{num_steps} = 20$). Then it will take iterations for the learner to go through all sentences once. Where N is the size of the list of words, b is batch size, and h is size of each sentence. So, the number of iterators is 1548
- Each batch data is read from train dataset of size 600, and shape of $[30 \times 20]$

First we start an interactive session:

In [9]:

```
session = tf.InteractiveSession()
```

```
/home/jupyterlab/conda/envs/python/lib/python3.6/site-
packages/tensorflow/python/client/session.py:1711:
UserWarning: An interactive session is already active. This
can cause out-of-memory errors in some cases. You must
explicitly call `InteractiveSession.close()` to release
```

```
resources held by the other session(s).
warnings.warn('An interactive session is already active.
This can '
```

```
In [10]:
```

```
# Reads the data and separates it into training data,
validation data and testing data
raw_data = reader.ptb_raw_data(data_dir)
train_data, valid_data, test_data, vocab, word_to_id =
raw_data
```

```
In [11]:
```

```
len(train_data)
```

```
Out[11]:
```

```
929589
```

```
In [12]:
```

```
def id_to_word(id_list):
    line = []
    for w in id_list:
        for word, wid in word_to_id.items():
            if wid == w:
                line.append(word)
    return line
```

```
print(id_to_word(train_data[0:100]))
```

```
['aer', 'banknote', 'berlitz', 'calloway', 'centrust',
'cluet', 'fromstein', 'gitano', 'guterman', 'hydro-
quebec', 'ipo', 'kia', 'memotec', 'mlx', 'nahb', 'punts',
'rake', 'regatta', 'rubens', 'sim', 'snack-food',
'ssangyong', 'swapo', 'wachter', '<eos>', 'pierre',
'<unk>', 'N', 'years', 'old', 'will', 'join', 'the',
'board', 'as', 'a', 'nonexecutive', 'director', 'nov.',
'N', '<eos>', 'mr.', '<unk>', 'is', 'chairman', 'of',
'<unk>', 'n.v.', 'the', 'dutch', 'publishing', 'group',
'<eos>', 'rudolph', '<unk>', 'N', 'years', 'old', 'and',
'former', 'chairman', 'of', 'consolidated', 'gold',
'fields', 'plc', 'was', 'named', 'a', 'nonexecutive',
'director', 'of', 'this', 'british', 'industrial',
'conglomerate', '<eos>', 'a', 'form', 'of', 'asbestos',
'once', 'used', 'to', 'make', 'kent', 'cigarette',
'filters', 'has', 'caused', 'a', 'high', 'percentage',
'of', 'cancer', 'deaths', 'among', 'a', 'group', 'of']
```

Lets just read one mini-batch now and feed our network:

```
In [13]:
```

```
itera = reader.ptb_iterator(train_data, batch_size,
```



```

num_steps)
first_touple = itera.__next__()
x = first_touple[0]
y = first_touple[1]
In [14]:
x.shape
Out[14]:
(60, 20)

```

Lets look at 3 sentences of our input x:

```

In [15]:
x[0:3]
Out[15]:
array([[9970, 9971, 9972, 9974, 9975, 9976, 9980, 9981,
        9982, 9983, 9984,
        9986, 9987, 9988, 9989, 9991, 9992, 9993, 9994,
        9995],
       [ 901,   33, 3361,    8, 1279,  437,  597,    6,
        261, 4276, 1089,
         8, 2836,    2,  269,    4, 5526,  241,   13,
        2420],
       [2654,    6,  334, 2886,    4,    1,  233,  711,
        834,   11,  130,
        123,    7,  514,    2,   63,   10,  514,    8,
        605]],
      dtype=int32)

```

we define 2 place holders to feed them with mini-batches, that is x and y:

```

In [16]:
_input_data = tf.placeholder(tf.int32, [batch_size,
num_steps]) #[30#20]
_targets = tf.placeholder(tf.int32, [batch_size,
num_steps]) #[30#20]

```

Lets define a dictionary, and use it later to feed the placeholders with our first mini-batch:

```

In [17]:
feed_dict = {_input_data:x, _targets:y}

```

For example, we can use it to feed _input_data:

```

In [18]:
session.run(_input_data, feed_dict)
Out[18]:
array([[9970, 9971, 9972, ..., 9993, 9994, 9995],

```

```

    [ 901,    33, 3361, ..., 241,    13, 2420],
    [2654,     6,  334, ..., 514,     8,  605],
    ...,
    [7831,    36, 1678, ...,    4, 4558,  157],
    [  59, 2070, 2433, ...,  400,    1, 1173],
    [2097,     3,    2, ..., 2043,    23,    1]],
dtype=int32)

```

In this step, we create the stacked LSTM, which is a 2 layer LSTM network:

```

In [19]:
lstm_cell_l1 = tf.contrib.rnn.BasicLSTMCell(hidden_size_l1,
forget_bias=0.0)
lstm_cell_l2 = tf.contrib.rnn.BasicLSTMCell(hidden_size_l2,
forget_bias=0.0)
stacked_lstm = tf.contrib.rnn.MultiRNNCell([lstm_cell_l1,
lstm_cell_l2])

```

Also, we initialize the states of the network:

_initial_state

For each LCTM, there are 2 state matrices, c_state and m_state. c_state and m_state represent "Memory State" and "Cell State". Each hidden layer, has a vector of size 30, which keeps the states. so, for 200 hidden units in each LSTM, we have a matrix of size [30x200]

```

In [20]:
_initial_state = stacked_lstm.zero_state(batch_size,
tf.float32)
_initial_state
Out[20]:
(LSTMStateTuple(c=<tf.Tensor 'MultiRNNCellZeroState/
BasicLSTMCellZeroState/zeros:0' shape=(60, 256)
dtype=float32>, h=<tf.Tensor 'MultiRNNCellZeroState/
BasicLSTMCellZeroState/zeros_1:0' shape=(60, 256)
dtype=float32>),
LSTMStateTuple(c=<tf.Tensor 'MultiRNNCellZeroState/
BasicLSTMCellZeroState_1/zeros:0' shape=(60, 128)
dtype=float32>, h=<tf.Tensor 'MultiRNNCellZeroState/
BasicLSTMCellZeroState_1/zeros_1:0' shape=(60, 128)
dtype=float32>))

```

Lets look at the states, though they are all zero for now:

```

In [21]:
session.run(_initial_state, feed_dict)

```

Out[21]:

```
(LSTMStateTuple(c=array([[0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          ...,
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.]]), dtype=float32),
h=array([[0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         ...,
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.]]), dtype=float32)),
LSTMStateTuple(c=array([[0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          ...,
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.],
                          [0., 0., 0., ..., 0., 0., 0.]]), dtype=float32),
h=array([[0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         ...,
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.],
         [0., 0., 0., ..., 0., 0., 0.]]), dtype=float32)))
```

Embeddings

We have to convert the words in our dataset to vectors of numbers. The traditional approach is to use one-hot encoding method that is usually used for converting categorical values to numerical values. However, One-hot encoded vectors are high-dimensional, sparse and in a big dataset, computationally inefficient. So, we use word2vec approach. It is, in fact, a layer in our LSTM network, where the word IDs will be represented as a dense representation before feeding to the LSTM.

The embedded vectors also get updated during the training process of the deep neural network. We create the embeddings for our input data. **embedding_vocab** is matrix of [10000x200] for all 10000 unique words.

In [22]:

```
embedding_vocab = tf.get_variable("embedding_vocab",
                                  [vocab_size, embedding_vector_size]) #[10000x200]
```

Lets initialize the `embedding_words` variable with random values.

In [23]:

```
session.run(tf.global_variables_initializer())
session.run(embedding_vocab)
```

Out[23]:

```
array([[ -0.01415229,  0.0020615 ,  0.00483348, ...,
        -0.02174182,
           -0.01958222, -0.0019541 ],
       [ 0.01009297, -0.0221437 ,  0.01230779, ...,
        -0.00450121,
           0.02269965, -0.01104087],
       [ 0.00581276, -0.00230083,  0.0093035 , ...,
        -0.00537598,
          -0.00462674,  0.00240099],
       ...,
       [ -0.02328045, -0.01938979, -0.02192219, ...,
        0.00431917,
          -0.00011383, -0.00199659],
       [ 0.00269953, -0.02154553,  0.01704296, ...,
        -0.00894606,
          -0.0072019 , -0.00217836],
       [ 0.00376112,  0.00271794, -0.009487 , ...,
        0.00215217,
          -0.01371387,  0.01537636]], dtype=float32)
```

embedding_lookup() finds the embedded values for our batch of 30x20 words. It goes to each row of `input_data`, and for each word in the row/sentence, finds the correspond vector in `embedding_dic`.

It creates a [30x20x200] tensor, so, the first element of **inputs** (the first sentence), is a matrix of 20x200, which each row of it, is vector representing a word in the sentence.

In [24]:

```
# Define where to get the data for our embeddings from
inputs = tf.nn.embedding_lookup(embedding_vocab,
    _input_data) #shape=(30, 20, 200)
inputs
```

Out[24]:

```
<tf.Tensor 'embedding_lookup:0' shape=(60, 20, 200)
dtype=float32>
```

In [25]:

```
session.run(inputs[0], feed_dict)
```

Out[25]:

```
array([[ -0.0064902 , -0.01207784, -0.00260266, ...,
```

```

0.0206253 ,
    -0.00664146,  0.01388235],
    [ 0.00208091, -0.01093761, -0.02060855, ...,
-0.0176512 ,
    0.01555531,  0.02231719],
    [ 0.00229493,  0.00974536,  0.00411461, ...,
-0.02234322,
    0.00752835, -0.01729599],
    ...,
    [-0.00346862, -0.01429997,  0.01049082, ...,
-0.01083067,
    0.01560483, -0.004133  ],
    [-0.01236486,  0.01851235, -0.00048184, ...,
-0.0117697 ,
    0.01201736,  0.01487473],
    [-0.01379531, -0.00858414,  0.00865652, ...,
-0.01257758,
    0.00151328, -0.00916983]], dtype=float32)

```

Constructing Recurrent Neural Networks

tf.nn.dynamic_rnn() creates a recurrent neural network using **stacked_lstm**.

The input should be a Tensor of shape: [batch_size, max_time, embedding_vector_size], in our case it would be (30, 20, 200)

This method, returns a pair (outputs, new_state) where:

- **outputs:** is a length T list of outputs (one for each input), or a nested tuple of such elements.
- **new_state:** is the final state.

In [26]:

```

outputs, new_state = tf.nn.dynamic_rnn(stacked_lstm,
inputs, initial_state=_initial_state)

```

so, lets look at the outputs. The output of the stackedLSTM comes from 200 hidden_layer, and in each time step(=20), one of them get activated. we use the linear activation to map the 200 hidden layer to a [?x10 matrix]

In [27]:

```

outputs

```

Out[27]:

```

<tf.Tensor 'rnn/transpose_1:0' shape=(60, 20, 128)
dtype=float32>

```

In [28]:

```

session.run(tf.global_variables_initializer())

```

```

session.run(outputs[0], feed_dict)
Out[28]:
array([[ -1.19189201e-04, -2.21901515e-04,
        -2.66724441e-04, ...,
         1.79185343e-04, -1.39329859e-04, -5.54598053e-04],
       [-2.98619503e-04, -6.10693125e-04,
        -4.65749530e-04, ...,
         -6.71920279e-05,  2.96846847e-04, -2.63071743e-05],
       [-5.01525719e-05, -9.99751966e-04,
        1.21005702e-04, ...,
         -5.49632859e-05,  6.09378854e-04,  6.72306051e-04],
       ...,
       [-4.19432763e-04, -7.15446193e-04,
        -1.88918988e-04, ...,
         -1.87663233e-03, -2.64889008e-04,  7.57927482e-04],
       [-7.73191336e-04, -1.10169302e-03,
        -1.68767765e-05, ...,
         -1.64000061e-03, -1.43042445e-04,  3.22060572e-04],
       [-4.49966057e-04, -9.23670596e-04,
        1.55205416e-04, ...,
         -1.64465094e-03,  1.15167524e-04,
        1.92213876e-04]], dtype=float32)

```

we need to flatten the outputs to be able to connect it softmax layer. Lets reshape the output tensor from [30 x 20 x 200] to [600 x 200].

Notice: Imagine our output is 3-d tensor as following (of course each `sen_x_word_y` is a an embedded vector by itself):

- sentence 1: [[sen1word1], [sen1word2], [sen1word3], ..., [sen1word20]]
- sentence 2: [[sen2word1], [sen2word2], [sen2word3], ..., [sen2word20]]
- sentence 3: [[sen3word1], [sen3word2], [sen3word3], ..., [sen3word20]]
- ...
- sentence 30: [[sen30word1], [sen30word2], [sen30word3], ..., [sen30word20]]

Now, the flatten would convert this 3-dim tensor to:

```

[ [sen1word1], [sen1word2], [sen1word3], ..., [sen1word20], [sen2word1], [sen2word2],
[sen2word3], ..., [sen2word20], ..., [sen30word20] ]

```

In [29]:

```
output = tf.reshape(outputs, [-1, hidden_size_l2])
```

output

Out[29]:

```
<tf.Tensor 'Reshape:0' shape=(1200, 128) dtype=float32>
```

logistic unit

Now, we create a logistic unit to return the probability of the output word in our vocabulary with 1000 words.

In [30]:

```
softmax_w = tf.get_variable("softmax_w", [hidden_size_12,
vocab_size]) #[200x1000]
softmax_b = tf.get_variable("softmax_b", [vocab_size])
#[1x1000]
logits = tf.matmul(output, softmax_w) + softmax_b
prob = tf.nn.softmax(logits)
```

Lets look at the probability of observing words for t=0 to t=20:

In [31]:

```
session.run(tf.global_variables_initializer())
output_words_prob = session.run(prob, feed_dict)
print("shape of the output: ", output_words_prob.shape)
print("The probability of observing words in t=0 to t=20",
output_words_prob[0:20])
```

```
shape of the output: (1200, 10000)
The probability of observing words in t=0 to t=20
[[1.00266683e-04 1.01174359e-04 9.92664136e-05 ...
9.82952843e-05
 9.97438692e-05 1.00532918e-04]
 [1.00270976e-04 1.01170968e-04 9.92712085e-05 ...
9.83011487e-05
 9.97425086e-05 1.00526580e-04]
 [1.00268669e-04 1.01171208e-04 9.92611822e-05 ...
9.83010177e-05
 9.97433453e-05 1.00523663e-04]
 ...
 [1.00266407e-04 1.01175487e-04 9.92567730e-05 ...
9.82916899e-05
 9.97421885e-05 1.00528479e-04]
 [1.00250683e-04 1.01166479e-04 9.92604619e-05 ...
9.82965430e-05
 9.97485040e-05 1.00528669e-04]
 [1.00248260e-04 1.01167454e-04 9.92575297e-05 ...
9.82902857e-05
 9.97535753e-05 1.00524267e-04]]
```

Prediction

What is the word correspond to the probability output? Lets use the maximum probability:

In [32]:

```
np.argmax(output_words_prob[0:20], axis=1)
```

Out[32]:

```
array([2713, 2713, 3913, 3913, 7784, 6848, 6848, 7784,
       616,  616, 5706,
        5706, 9432, 9432, 9432, 2073, 2073, 2073, 9432,
       2073])
```

So, what is the ground truth for the first word of first sentence?

In [33]:

```
y[0]
```

Out[33]:

```
array([9971, 9972, 9974, 9975, 9976, 9980, 9981, 9982,
       9983, 9984, 9986,
        9987, 9988, 9989, 9991, 9992, 9993, 9994, 9995,
       9996], dtype=int32)
```

Also, you can get it from target tensor, if you want to find the embedding vector:

In [34]:

```
targ = session.run(_targets, feed_dict)
```

```
targ[0]
```

Out[34]:

```
array([9971, 9972, 9974, 9975, 9976, 9980, 9981, 9982,
       9983, 9984, 9986,
        9987, 9988, 9989, 9991, 9992, 9993, 9994, 9995,
       9996], dtype=int32)
```

How similar the predicted words are to the target words?

Objective function

Now we have to define our objective function, to calculate the similarity of predicted values to ground truth, and then, penalize the model with the error. Our objective is to minimize loss function, that is, to minimize the average negative log probability of the target words:

This function is already implemented and available in TensorFlow through **sequence_loss_by_example**. It calculates the weighted cross-entropy loss for **logits** and the **target** sequence.

The arguments of this function are:

- logits: List of 2D Tensors of shape [batch_size x num_decoder_symbols].
- targets: List of 1D batch-sized int32 Tensors of the same length as logits.
- weights: List of 1D batch-sized float-Tensors of the same length as logits.

In [35]:

```
loss =
tf.contrib.legacy_seq2seq.sequence_loss_by_example([logits]
, [tf.reshape(_targets, [-1])],[tf.ones([batch_size *
num_steps])])
```

loss is a 1D batch-sized float Tensor [600x1]: The log-perplexity for each sequence. Lets look at the first 10 values of loss:

In [36]:

```
session.run(loss, feed_dict)[:10]
```

Out[36]:

```
array([ 9.205525 ,  9.2144985,  9.221913 ,  9.202395 ,
        9.227202 ,  9.210929 ,
        9.209447 ,  9.1944895,  9.210582 ,  9.210322 ],
      dtype=float32)
```

Now, we define loss as average of the losses:

In [37]:

```
cost = tf.reduce_sum(loss) / batch_size
session.run(tf.global_variables_initializer())
session.run(cost, feed_dict)
```

Out[37]:

```
184.25804
```

Training

To do training for our network, we have to take the following steps:

1. Define the optimizer.
2. Extract variables that are trainable.
3. Calculate the gradients based on the loss function.
4. Apply the optimizer to the variables/gradients tuple.

1. Define Optimizer

GradientDescentOptimizer constructs a new gradient descent optimizer. Later, we use constructed **optimizer** to compute gradients for a loss and apply gradients to variables.

In [38]:

```
# Create a variable for the learning rate
lr = tf.Variable(0.0, trainable=False)
# Create the gradient descent optimizer with our learning
```

```
rate
optimizer = tf.train.GradientDescentOptimizer(lr)
```

2. Trainable Variables

Defining a variable, if you passed *trainable=True*, the variable constructor automatically adds new variables to the graph collection **GraphKeys.TRAINABLE_VARIABLES**. Now, using *tf.trainable_variables()* you can get all variables created with **trainable=True**.

In [39]:

```
# Get all TensorFlow variables marked as "trainable" (i.e.
all of them except _lr, which we just created)
tvars = tf.trainable_variables()
```

tvars

Out[39]:

```
[<tf.Variable 'embedding_vocab:0' shape=(10000, 200)
dtype=float32_ref>,
 <tf.Variable 'rnn/multi_rnn_cell/cell_0/basic_lstm_cell/
kernel:0' shape=(456, 1024) dtype=float32_ref>,
 <tf.Variable 'rnn/multi_rnn_cell/cell_0/basic_lstm_cell/
bias:0' shape=(1024,) dtype=float32_ref>,
 <tf.Variable 'rnn/multi_rnn_cell/cell_1/basic_lstm_cell/
kernel:0' shape=(384, 512) dtype=float32_ref>,
 <tf.Variable 'rnn/multi_rnn_cell/cell_1/basic_lstm_cell/
bias:0' shape=(512,) dtype=float32_ref>,
 <tf.Variable 'softmax_w:0' shape=(128, 10000)
dtype=float32_ref>,
 <tf.Variable 'softmax_b:0' shape=(10000,)
dtype=float32_ref>]
```

Note: we can find the name and scope of all variables:

In [40]:

```
[v.name for v in tvars]
```

Out[40]:

```
['embedding_vocab:0',
 'rnn/multi_rnn_cell/cell_0/basic_lstm_cell/kernel:0',
 'rnn/multi_rnn_cell/cell_0/basic_lstm_cell/bias:0',
 'rnn/multi_rnn_cell/cell_1/basic_lstm_cell/kernel:0',
 'rnn/multi_rnn_cell/cell_1/basic_lstm_cell/bias:0',
 'softmax_w:0',
 'softmax_b:0']
```

3. Calculate the gradients based on the loss function

Gradient

: The gradient of a function is the slope of its derivative (line), or in other words, the rate of change of a function. It's a vector (a direction to move) that points in the direction of greatest increase of the function, and calculated by the **derivative** operation.

First lets recall the gradient function using an toy example:

In [41]:

```
var_x = tf.placeholder(tf.float32)
var_y = tf.placeholder(tf.float32)
func_test = 2.0 * var_x * var_x + 3.0 * var_x * var_y
session.run(tf.global_variables_initializer())
session.run(func_test, {var_x:1.0,var_y:2.0})
```

Out[41]:

8.0

The **tf.gradients()** function allows you to compute the symbolic gradient of one tensor with respect to one or more other tensors—including variables. **tf.gradients(func, xs)** constructs symbolic partial derivatives of sum of **func** w.r.t. x in **xs**.

Now, lets look at the derivative w.r.t. **var_x**:

In [42]:

```
var_grad = tf.gradients(func_test, [var_x])
session.run(var_grad, {var_x:1.0,var_y:2.0})
```

Out[42]:

[10.0]

the derivative w.r.t. **var_y**:

In [43]:

```
var_grad = tf.gradients(func_test, [var_y])
session.run(var_grad, {var_x:1.0, var_y:2.0})
```

Out[43]:

[3.0]

Now, we can look at gradients w.r.t all variables:

In [44]:

```
tf.gradients(cost, tvars)
```

Out[44]:

```
[<tensorflow.python.framework.ops.IndexedSlices at
0x7f586c0e1d68>,
 <tf.Tensor 'gradients_2/rnn/while/rnn/multi_rnn_cell/
cell_0/basic_lstm_cell/MatMul/Enter_grad/b_acc_3:0'
```

```

shape=(456, 1024) dtype=float32>,
<tf.Tensor 'gradients_2/rnn/while/rnn/multi_rnn_cell/
cell_0/basic_lstm_cell/BiasAdd/Enter_grad/b_acc_3:0'
shape=(1024,) dtype=float32>,
<tf.Tensor 'gradients_2/rnn/while/rnn/multi_rnn_cell/
cell_1/basic_lstm_cell/MatMul/Enter_grad/b_acc_3:0'
shape=(384, 512) dtype=float32>,
<tf.Tensor 'gradients_2/rnn/while/rnn/multi_rnn_cell/
cell_1/basic_lstm_cell/BiasAdd/Enter_grad/b_acc_3:0'
shape=(512,) dtype=float32>,
<tf.Tensor 'gradients_2/MatMul_grad/MatMul_1:0'
shape=(128, 10000) dtype=float32>,
<tf.Tensor 'gradients_2/add_grad/Reshape_1:0'
shape=(10000,) dtype=float32>]

```

In [45]:

```

grad_t_list = tf.gradients(cost, tvars)
#sess.run(grad_t_list, feed_dict)

```

now, we have a list of tensors, *t-list*. We can use it to find clipped tensors.

clip_by_global_norm clips values of multiple tensors by the ratio of the sum of their norms.

clip_by_global_norm get *t-list* as input and returns 2 things:

- a list of clipped tensors, so called *list_clipped*
- the global norm (*global_norm*) of all tensors in *t_list*

In [46]:

```

# Define the gradient clipping threshold
grads, _ = tf.clip_by_global_norm(grad_t_list,
max_grad_norm)
grads

```

Out[46]:

```

[<tensorflow.python.framework.ops.IndexedSlices at
0x7f586c083400>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_1:0'
shape=(456, 1024) dtype=float32>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_2:0'
shape=(1024,) dtype=float32>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_3:0'
shape=(384, 512) dtype=float32>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_4:0'
shape=(512,) dtype=float32>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_5:0'
shape=(128, 10000) dtype=float32>,
<tf.Tensor 'clip_by_global_norm/clip_by_global_norm/_6:0'
shape=(10000,) dtype=float32>]

```

In [47]:

```
session.run(grads, feed_dict)
```

```
Out[47]:
```

```
[IndexedSlicesValue(values=array([[ -2.61569548e-05,
 6.16112402e-06, -2.66044071e-06, ...,
      1.74453453e-05, -1.04860119e-05,
 5.10769542e-06],
      [-1.77246911e-05,  1.21485164e-05,
-1.38116206e-06, ...,
      1.11905920e-05, -5.60408034e-06,
 7.72732164e-06],
      [-8.68449661e-06,  1.11302970e-05,
 6.31111607e-06, ...,
      8.65887728e-07,  1.27591832e-07,
 6.00685280e-06],
      ...,
      [ 3.11789017e-06, -6.15000044e-06,
 1.26706154e-05, ...,
      -1.25420165e-05,  4.80279095e-06,
 5.42298221e-06],
      [ 1.50459653e-06, -8.42342706e-06,
 7.66077301e-06, ...,
      -1.07366031e-05, -4.16762259e-06,
 7.37164510e-06],
      [-2.84778866e-06, -3.81677637e-06,
 7.80785467e-06, ...,
      -1.55558564e-06, -9.58248734e-07,
 5.77966830e-06]], dtype=float32), indices=array([9970,
 9971, 9972, ..., 2043,   23,    1], dtype=int32),
dense_shape=array([10000,   200], dtype=int32)),
array([[ 1.6066767e-08,  3.5067192e-08,
-1.1746049e-10, ...,
      5.8436541e-08, -2.3235570e-08, -3.1973450e-08],
      [-5.5296660e-08, -3.3499763e-08,
 3.5328263e-08, ...,
      -2.5721681e-08,  2.9637341e-08,  4.8090200e-08],
      [-8.4972394e-09, -2.5726406e-08,
 6.3189923e-08, ...,
      -2.8150120e-08, -1.6127135e-09,  4.9952789e-08],
      ...,
      [ 5.2933906e-09, -2.5581681e-09,
 2.5126173e-10, ...,
      1.1959569e-08, -6.4872872e-09,  1.1115321e-08],
      [-2.3990161e-09,  5.1252247e-10,
-2.2163025e-09, ...,
      -4.3951234e-09, -2.0027489e-08, -5.1065313e-10],
      [-6.4774719e-09, -3.9391477e-09,
 4.1060377e-09, ...,
```

```

        4.9047104e-09, -5.1072515e-09,  1.4783470e-08]],
dtype=float32),
    array([ 1.1906226e-06,  1.2429670e-06,
-3.3754725e-07, ...,
        3.4770521e-06,  2.6483815e-06,  1.4185525e-06],
dtype=float32),
    array([[ 2.82152257e-09,  2.86482749e-09,
-3.88555987e-09, ...,
        -4.84896434e-09,  1.49613233e-09,
-3.22891602e-09],
        [ 3.33813843e-09,  8.23170598e-09,
4.96051289e-09, ...,
        -7.32058947e-09, -7.41958850e-09,
3.94051280e-09],
        [-1.12334786e-09,  3.23000893e-09,
1.58423274e-09, ...,
        -2.84767854e-09, -1.09925735e-10,
-2.44763632e-09],
        ...,
        [-7.39758144e-10,  3.41527673e-09,
-1.07006737e-09, ...,
        -5.25751942e-10, -1.61392688e-09,
2.94533647e-10],
        [ 1.31301436e-09, -2.85678481e-10,
-6.59250876e-10, ...,
        -8.46696158e-10, -7.34119121e-09,
-1.56427304e-09],
        [-1.62541258e-10,  9.05795772e-10,
2.15344617e-10, ...,
        -2.96144687e-09, -3.10188675e-09,
9.89087834e-09]], dtype=float32),
    array([ 4.69054112e-06,  4.68102058e-07, -9.92032938e-07,
5.60780336e-07,
        -2.59410217e-06, -2.22886888e-06,  2.45557612e-06,
-8.80732955e-07,
        -4.70148223e-08,  4.76234709e-06, -3.70097268e-07,
-2.24516452e-06,
        -3.53951759e-06,  7.30050715e-07,  3.52373775e-07,
2.82506198e-06,
        1.71158308e-06,  2.37686118e-06, -1.92323273e-06,
-2.41163934e-06,
        1.60618094e-06,  9.53752647e-07,  4.75565230e-06,
5.22564619e-07,
        3.44931459e-06,  4.31589683e-07,  6.04613479e-06,
-2.25653025e-06,
        6.04737852e-06, -8.20218611e-06,  1.56185456e-06,
3.34413653e-06,

```

1.52187476e-06, -7.72086150e-06, 2.58739783e-06,
-9.71262352e-07,
7.35883532e-06, 1.55628629e-06, -3.24571907e-08,
-3.44389582e-06,
2.01040643e-06, -9.10827566e-07, 1.17336839e-07,
8.35887658e-06,
-2.45897195e-06, -2.53624103e-06, 2.34787217e-06,
-2.08731808e-06,
1.73975025e-07, 7.02792431e-06, 5.80364303e-06,
-8.19647448e-06,
-4.06862142e-07, 2.30448222e-06, -3.65284473e-06,
7.44432145e-06,
5.23486824e-06, -1.98140970e-06, -2.61483933e-07,
6.83807571e-07,
-7.05898856e-06, 3.94286144e-06, 4.43415047e-07,
6.91718355e-07,
6.51707410e-07, -6.59503826e-07, 1.59254427e-07,
-5.82457278e-06,
1.97188115e-06, -4.30781893e-06, 6.94200196e-07,
1.51332915e-06,
1.62943354e-06, -1.78604489e-06, 4.05077515e-07,
1.47902506e-06,
-4.87039188e-06, -1.63671677e-06, -4.06812296e-06,
1.08095878e-06,
3.04282685e-06, 1.11368149e-06, 4.61136506e-06,
5.18808633e-07,
-1.17497575e-06, -2.02238107e-06, 1.82892177e-06,
1.90911760e-06,
-3.12053521e-06, -3.55679867e-06, -2.95036489e-06,
6.15979798e-06,
5.11595317e-08, 1.79633309e-06, 8.06270066e-07,
1.15697208e-06,
-2.28780527e-06, 1.64003427e-06, -6.33194929e-07,
6.09375866e-06,
-6.34333583e-06, 3.52051103e-07, -5.96862446e-06,
-1.47543915e-06,
-7.09807682e-06, -9.97329153e-07, 3.48767890e-06,
1.11682402e-06,
-5.72341469e-06, 4.75494153e-06, 5.31957357e-06,
-1.00912462e-06,
-4.58540671e-06, -9.97827556e-07, 2.29915736e-06,
-1.90672267e-06,
4.84219072e-07, -9.30472538e-07, 4.36010032e-06,
5.73212947e-06,
1.12038651e-06, -6.01142756e-06, -1.80933591e-06,
8.50413358e-07,
-3.40388624e-06, -4.88233718e-06, -2.65788958e-06,

-4.23499250e-06,
-1.67782698e-02, 1.40324084e-03, -2.05439292e-02,
-1.56230200e-02,
-6.46257307e-03, 2.13513356e-02, 2.61986221e-04,
-6.80318754e-03,
8.11728928e-03, -2.36294698e-02, 1.34925721e-02,
-6.54587476e-03,
-2.40131449e-02, -2.80507677e-03, -1.14419786e-02,
-2.15098262e-02,
6.07113261e-03, -8.78779124e-03, -1.12909470e-02,
-1.30228950e-02,
6.79672044e-03, -7.87025201e-04, -1.54209752e-02,
-1.13330770e-03,
6.75834669e-03, -2.03032754e-02, -1.71342976e-02,
-1.70419086e-02,
-6.20212522e-04, -4.47902940e-02, -1.23471338e-02,
-1.36425393e-02,
2.72955038e-02, 1.62695069e-02, 2.00330019e-02,
-1.94328446e-02,
1.48442900e-02, 7.48993829e-03, 1.59891397e-02,
1.93818249e-02,
1.45618357e-02, -4.68219118e-03, -1.07645923e-02,
-2.57405415e-02,
-1.12656774e-02, 5.06304624e-03, 2.43847836e-02,
2.09896471e-02,
-3.65569303e-03, 3.23981717e-02, -2.11831108e-02,
-2.16212459e-02,
1.83932623e-03, -9.95439477e-03, 1.38720982e-02,
3.85212749e-02,
-1.83281209e-03, 2.29631439e-02, -6.64313324e-03,
7.27478997e-04,
-2.25201733e-02, 2.22989451e-02, 2.39708787e-03,
-5.05504990e-03,
8.49392451e-03, 8.86597577e-03, 1.13722973e-03,
1.74824353e-02,
1.91654935e-02, -1.33035919e-02, 1.20981587e-02,
9.90184862e-03,
-1.21057807e-02, 6.36008568e-03, -9.03539546e-03,
1.12846782e-02,
1.09304842e-02, 1.09546296e-02, -6.41993899e-03,
1.09641282e-02,
2.13901065e-02, 1.60594180e-04, -1.46217626e-02,
4.13014274e-03,
1.29602570e-02, 1.54162208e-02, 1.77017748e-02,
2.54471600e-02,
-2.89960224e-02, -2.29745246e-02, -5.75283915e-03,
2.03590151e-02,

3.25137414e-02, -1.80340596e-02, 1.40880691e-02,
3.62778665e-03,
-5.94495470e-03, -2.83740670e-03, 1.78599823e-02,
-1.38584413e-02,
1.76584721e-02, -7.64559535e-03, 2.13941969e-02,
-2.89646275e-02,
-6.42309431e-03, -1.53222857e-02, 1.37541126e-02,
6.65195240e-03,
-2.51375921e-02, -8.90096184e-03, -3.81129645e-02,
-1.56557150e-02,
-2.47294698e-02, 1.97428744e-03, 4.07801569e-03,
1.50190741e-02,
6.92302128e-03, 2.18006894e-02, 8.67008790e-03,
1.23183271e-02,
1.08832121e-02, -3.18964571e-02, 3.39374598e-03,
-8.71196575e-03,
1.36874262e-02, 3.04737571e-03, -2.95267683e-02,
3.68994884e-02,
5.75129548e-07, -2.63159109e-07, 1.30382750e-06,
-9.28824534e-07,
-3.41057671e-06, 2.60241961e-07, 1.64374831e-06,
-1.34782113e-06,
-1.17480141e-07, 3.49939455e-06, -4.55303962e-07,
-5.16705722e-06,
-4.99210728e-06, 2.11140582e-06, 8.88680688e-07,
3.36226231e-06,
1.55396435e-06, 3.71765282e-06, -1.67785709e-06,
-2.65824474e-06,
2.59022158e-07, -3.47943842e-06, 4.32638262e-06,
4.10687733e-07,
2.34934123e-06, 1.37174493e-06, 5.59340742e-06,
-4.85302007e-06,
4.22242874e-06, -6.43024032e-06, -8.51564891e-07,
3.74586648e-06,
2.31881899e-07, -8.15783642e-06, 3.52595930e-06,
-2.28803833e-06,
3.69596364e-06, 7.36391939e-07, -2.94863924e-07,
-4.43130421e-06,
-3.57236013e-06, -1.50326787e-06, 3.69040663e-07,
6.46886019e-06,
-1.46349834e-07, 3.83327887e-07, 1.44465071e-06,
-1.35943628e-06,
-2.25721536e-07, 5.16130012e-06, 4.72013789e-06,
-8.52984613e-06,
-3.62786665e-07, 2.52926520e-06, -2.52609971e-06,
3.78092000e-06,
1.33412561e-06, -3.63973595e-06, -2.62799290e-06,

-1.17601621e-06,
-4.63347988e-06, -8.29023350e-07, 1.44919602e-06,
1.68958547e-06,
1.84339240e-06, -7.98195344e-07, 6.55083738e-08,
-4.06115578e-06,
1.74780996e-06, -3.86603961e-06, -1.16461976e-07,
-2.94971869e-07,
3.06160246e-06, -2.27214878e-06, -3.80240590e-06,
3.58638931e-06,
-5.06503420e-06, -6.37957237e-06, -1.63210848e-06,
1.56711292e-06,
9.60996886e-07, -2.59937087e-06, 4.37669996e-06,
1.65900587e-06,
-1.83292957e-06, -2.03104423e-06, 6.09795791e-07,
-7.88546799e-07,
3.92503807e-11, -3.40877523e-06, -2.22923563e-06,
5.14046133e-06,
-1.46384551e-07, 6.95680569e-07, -3.55125081e-07,
2.48304772e-08,
-7.78362676e-07, 8.15811120e-07, 2.60363322e-06,
5.00697570e-06,
-5.69264103e-06, -3.05187609e-06, -4.34491130e-06,
2.75907291e-06,
-4.87445868e-06, 1.36171684e-06, 1.29883415e-06,
5.83035444e-07,
-4.74623994e-06, 5.28778082e-06, 8.77900584e-06,
-1.58754119e-06,
6.65091534e-07, 6.23058497e-07, 1.34108427e-06,
-3.54326380e-06,
-1.36215522e-06, -2.73867363e-06, 4.22087123e-06,
7.85362317e-06,
-9.16213821e-07, -4.36956861e-06, -3.91406019e-07,
-2.39307838e-07,
-2.05247170e-06, -1.72826014e-06, -2.32609659e-06,
-3.46989736e-06,
4.69911993e-06, 4.65100896e-07, -9.89996124e-07,
5.53963957e-07,
-2.58802129e-06, -2.22813560e-06, 2.46172658e-06,
-8.85761324e-07,
-4.38511449e-08, 4.76039850e-06, -3.69381894e-07,
-2.23981351e-06,
-3.53272708e-06, 7.29258829e-07, 3.57510714e-07,
2.82647807e-06,
1.71171621e-06, 2.37360723e-06, -1.92001835e-06,
-2.41163229e-06,
1.60607181e-06, 9.57939278e-07, 4.75295792e-06,
5.15477552e-07,

3.45039143e-06, 4.33367518e-07, 6.03808530e-06,
-2.25976964e-06,
6.03929902e-06, -8.19775687e-06, 1.55941780e-06,
3.33800085e-06,
1.51998393e-06, -7.72514159e-06, 2.58542855e-06,
-9.77377795e-07,
7.36141328e-06, 1.55773660e-06, -3.70788200e-08,
-3.44189243e-06,
2.00991326e-06, -9.10174549e-07, 1.16632265e-07,
8.35962419e-06,
-2.46242143e-06, -2.53434609e-06, 2.35014431e-06,
-2.08532128e-06,
1.71112504e-07, 7.02001034e-06, 5.80326923e-06,
-8.19709294e-06,
-4.08654643e-07, 2.29968919e-06, -3.65223514e-06,
7.44271892e-06,
5.22901246e-06, -1.98165299e-06, -2.56819476e-07,
6.84374527e-07,
-7.05250568e-06, 3.93815253e-06, 4.43990729e-07,
6.89872991e-07,
6.53909751e-07, -6.59104558e-07, 1.69633154e-07,
-5.82187067e-06,
1.97146187e-06, -4.31157787e-06, 6.94333266e-07,
1.51266477e-06,
1.63146103e-06, -1.78288110e-06, 4.05664395e-07,
1.47696278e-06,
-4.86814497e-06, -1.63414097e-06, -4.06505706e-06,
1.07520077e-06,
3.04547734e-06, 1.11393342e-06, 4.61753280e-06,
5.20912920e-07,
-1.17195918e-06, -2.02326532e-06, 1.82895121e-06,
1.90987339e-06,
-3.13556939e-06, -3.55966358e-06, -2.94711640e-06,
6.15707904e-06,
4.83761369e-08, 1.79833864e-06, 8.09227515e-07,
1.15738396e-06,
-2.28903514e-06, 1.63878872e-06, -6.29824228e-07,
6.09275503e-06,
-6.34225535e-06, 3.59656525e-07, -5.96742984e-06,
-1.47614332e-06,
-7.09931146e-06, -9.99413942e-07, 3.49338825e-06,
1.11041788e-06,
-5.72565204e-06, 4.76047853e-06, 5.31512296e-06,
-1.00100249e-06,
-4.58489740e-06, -9.97561074e-07, 2.29854550e-06,
-1.90547996e-06,
4.84413590e-07, -9.27973133e-07, 4.36150003e-06,

```

5.73578745e-06,
    1.11430541e-06, -6.02246655e-06, -1.80911707e-06,
8.52821131e-07,
    -3.40345741e-06, -4.88256183e-06, -2.65874633e-06,
-4.23911843e-06],
    dtype=float32),
    array([[ -2.36724190e-05,  1.65698075e-04,
 1.51683998e-04, ...,
        -1.07361188e-07, -1.05886855e-07,
-1.05065588e-07],
        [-3.01451189e-04, -2.43887363e-04,
-2.50373443e-04, ...,
        5.98761744e-07,  5.90605680e-07,
 5.85834016e-07],
        [-1.49160687e-05, -6.38432393e-05,
 1.63396890e-05, ...,
        1.96839345e-09,  1.95364969e-09,
 1.93524574e-09],
        ...,
        [-6.59689031e-05,  9.05689667e-05,
-8.56744009e-07, ...,
        1.45642076e-07,  1.43621051e-07,
 1.42496731e-07],
        [-4.08373380e-05, -1.14960385e-04,
 1.56148293e-04, ...,
        -6.80299053e-08, -6.71384015e-08,
-6.65706850e-08],
        [ 1.13050970e-04,  6.82620157e-05,
 1.66142563e-04, ...,
        -1.12256231e-07, -1.10693804e-07,
-1.09840236e-07]], dtype=float32),
    array([-0.78133273, -1.0979931 , -0.98132986, ...,
 0.00202222,
        0.00199453,  0.00197843], dtype=float32)])

```

4. Apply the optimizer to the variables / gradients tuple.

In [48]:

```

# Create the training TensorFlow Operation through our
optimizer
train_op = optimizer.apply_gradients(zip(grads, tvars))

```

In [49]:

```

session.run(tf.global_variables_initializer())
session.run(train_op, feed_dict)

```

LSTM

We learned how the model is build step by step. Noe, let's then create a Class that represents our model. This class needs a few things:

- We have to create the model in accordance with our defined hyperparameters
- We have to create the placeholders for our input data and expected outputs (the real data)
- We have to create the LSTM cell structure and connect them with our RNN structure
- We have to create the word embeddings and point them to the input data
- We have to create the input structure for our RNN
- We have to instantiate our RNN model and retrieve the variable in which we should expect our outputs to appear
- We need to create a logistic structure to return the probability of our words
- We need to create the loss and cost functions for our optimizer to work, and then create the optimizer
- And finally, we need to create a training operation that can be run to actually train our model

In [50]:

```
hidden_size_l1
```

Out[50]:

```
256
```

In [51]:

```
class PTBModel(object):
```

```
    def __init__(self, action_type):
        #####
        # Setting parameters for ease of use #
        #####
        self.batch_size = batch_size
        self.num_steps = num_steps
        self.hidden_size_l1 = hidden_size_l1
        self.hidden_size_l2 = hidden_size_l2
        self.vocab_size = vocab_size
        self.embedding_vector_size = embedding_vector_size
```

```
#####
#####
```

```
        # Creating placeholders for our input data and
        expected outputs (target data) #
```

```
#####
#####
```

```
        self._input_data = tf.placeholder(tf.int32,
```

```

[batch_size, num_steps]) #[30#20]
    self._targets = tf.placeholder(tf.int32,
[batch_size, num_steps]) #[30#20]

#####
#####
    # Creating the LSTM cell structure and connect it
    with the RNN structure #

#####
#####
    # Create the LSTM unit.
    # This creates only the structure for the LSTM and
    has to be associated with a RNN unit still.
    # The argument n_hidden(size=200) of BasicLSTMCell
    is size of hidden layer, that is, the number of hidden
    units of the LSTM (inside A).
    # Size is the same as the size of our hidden layer,
    and no bias is added to the Forget Gate.
    # LSTM cell processes one word at a time and
    computes probabilities of the possible continuations of the
    sentence.
    lstm_cell_l1 =
tf.contrib.rnn.BasicLSTMCell(self.hidden_size_l1,
forget_bias=0.0)
    lstm_cell_l2 =
tf.contrib.rnn.BasicLSTMCell(self.hidden_size_l2,
forget_bias=0.0)

    # Unless you changed keep_prob, this won't actually
    execute -- this is a dropout wrapper for our LSTM unit
    # This is an optimization of the LSTM output, but
    is not needed at all
    if action_type == "is_training" and keep_prob < 1:
        lstm_cell_l1 =
tf.contrib.rnn.DropoutWrapper(lstm_cell_l1,
output_keep_prob=keep_prob)
        lstm_cell_l2 =
tf.contrib.rnn.DropoutWrapper(lstm_cell_l2,
output_keep_prob=keep_prob)

    # By taking in the LSTM cells as parameters, the
    MultiRNNCell function junctions the LSTM units to the RNN
    units.
    # RNN cell composed sequentially of multiple simple
    cells.

```

```

        stacked_lstm =
tf.contrib.rnn.MultiRNNCell([lstm_cell_l1, lstm_cell_l2])

        # Define the initial state, i.e., the model state
for the very first data point
        # It initialize the state of the LSTM memory. The
memory state of the network is initialized with a vector of
zeros and gets updated after reading each word.
        self._initial_state =
stacked_lstm.zero_state(batch_size, tf.float32)

#####
#####
        # Creating the word embeddings and pointing them to
the input data #

#####
#####
        with tf.device("/cpu:0"):
            # Create the embeddings for our input data.
Size is hidden size.
            embedding = tf.get_variable("embedding",
[vocab_size, self.embedding_vector_size]) #[10000x200]
            # Define where to get the data for our
embeddings from
            inputs = tf.nn.embedding_lookup(embedding,
self._input_data)

        # Unless you changed keep_prob, this won't actually
execute -- this is a dropout addition for our inputs
        # This is an optimization of the input processing
and is not needed at all
        if action_type == "is_training" and keep_prob < 1:
            inputs = tf.nn.dropout(inputs, keep_prob)

#####
        # Creating the input structure for our RNN #
#####
        # Input structure is 20x[30x200]
        # Considering each word is represented by a 200
dimensional vector, and we have 30 batchs, we create 30
word-vectors of size [30xx2000]
        # inputs = [tf.squeeze(input_, [1]) for input_ in
tf.split(1, num_steps, inputs)]
        # The input structure is fed from the embeddings,
which are filled in by the input data

```

```
        # Feeding a batch of b sentences to a RNN:
        # In step 1, first word of each of the b sentences
(in a batch) is input in parallel.
        # In step 2, second word of each of the b
sentences is input in parallel.
        # The parallelism is only for efficiency.
        # Each sentence in a batch is handled in parallel,
but the network sees one word of a sentence at a time and
does the computations accordingly.
        # All the computations involving the words of all
sentences in a batch at a given time step are done in
parallel.
```

```
#####
#####
```

```
        # Instantiating our RNN model and retrieving the
structure for returning the outputs and the state #
```

```
#####
#####
```

```
        outputs, state = tf.nn.dynamic_rnn(stacked_lstm,
inputs, initial_state=self._initial_state)
```

```
#####
#####
```

```
        # Creating a logistic unit to return the
probability of the output word #
```

```
#####
#####
```

```
        output = tf.reshape(outputs, [-1,
self.hidden_size_l2])
        softmax_w = tf.get_variable("softmax_w",
[self.hidden_size_l2, vocab_size]) #[200x1000]
        softmax_b = tf.get_variable("softmax_b",
[vocab_size]) #[1x1000]
        logits = tf.matmul(output, softmax_w) + softmax_b
        logits = tf.reshape(logits, [self.batch_size,
self.num_steps, vocab_size])
        prob = tf.nn.softmax(logits)
        out_words = tf.argmax(prob, axis=2)
        self._output_words = out_words
```

```
#####
#####
```



```

        # Defining the loss and cost functions for the
model's learning to work #

#####

        # Use the contrib sequence loss and average over
the batches
        loss = tf.contrib.seq2seq.sequence_loss(
            logits,
            self.targets,
            tf.ones([batch_size, num_steps],
dtype=tf.float32),
            average_across_timesteps=False,
            average_across_batch=True)

#         loss =
tf.contrib.legacy_seq2seq.sequence_loss_by_example([logits]
, [tf.reshape(self._targets, [-1])],
#
[tf.ones([batch_size * num_steps])])
        self._cost = tf.reduce_sum(loss)

        # Store the final state
self._final_state = state

        #Everything after this point is relevant only for
training
        if action_type != "is_training":
            return

#####
        # Creating the Training Operation for our Model #
#####
        # Create a variable for the learning rate
self._lr = tf.Variable(0.0, trainable=False)
        # Get all TensorFlow variables marked as
"trainable" (i.e. all of them except _lr, which we just
created)
        tvars = tf.trainable_variables()
        # Define the gradient clipping threshold
grads, _ =
tf.clip_by_global_norm(tf.gradients(self._cost, tvars),
max_grad_norm)
        # Create the gradient descent optimizer with our
learning rate

```

```

        optimizer =
tf.train.GradientDescentOptimizer(self.lr)
        # Create the training TensorFlow Operation through
our optimizer
        self._train_op =
optimizer.apply_gradients(zip(grads, tvars))

    # Helper functions for our LSTM RNN class

    # Assign the learning rate for this model
    def assign_lr(self, session, lr_value):
        session.run(tf.assign(self.lr, lr_value))

    # Returns the input data for this model at a point in
time
    @property
    def input_data(self):
        return self._input_data

    # Returns the targets for this model at a point in time
    @property
    def targets(self):
        return self._targets

    # Returns the initial state for this model
    @property
    def initial_state(self):
        return self._initial_state

    # Returns the defined Cost
    @property
    def cost(self):
        return self._cost

    # Returns the final state for this model
    @property
    def final_state(self):
        return self._final_state

    # Returns the final output words for this model
    @property
    def final_output_words(self):
        return self._output_words

    # Returns the current learning rate for this model

```

```

@property
def lr(self):
    return self._lr

# Returns the training operation defined for this model
@property
def train_op(self):
    return self._train_op

```

With that, the actual structure of our Recurrent Neural Network with Long Short-Term Memory is finished. What remains for us to do is to actually create the methods to run through time -- that is, the `run_epoch` method to be run at each epoch and a main script which ties all of this together.

What our `run_epoch` method should do is take our input data and feed it to the relevant operations. This will return at the very least the current result for the cost function.

In [52]:

```

#####
#####
####
# run_one_epoch takes as parameters the current session,
the model instance, the data to be fed, and the operation
to be run #
#####
#####
####
def run_one_epoch(session, m, data, eval_op,
verbose=False):

    #Define the epoch size based on the length of the data,
batch size and the number of steps
    epoch_size = ((len(data) // m.batch_size) - 1) //
m.num_steps
    start_time = time.time()
    costs = 0.0
    iters = 0

    state = session.run(m.initial_state)

    #For each step and data point
    for step, (x, y) in enumerate(reader.ptb_iterator(data,
m.batch_size, m.num_steps)):

        #Evaluate and return cost, state by running cost,
final_state and the function passed as parameter
        cost, state, out_words, _ = session.run([m.cost,

```

```

m.final_state, m.final_output_words, eval_op],
                                {m.input_data: x,
                                 m.targets: y,
                                 m.initial_state:
state})

    #Add returned cost to costs (which keeps track of
the total costs for this epoch)
    costs += cost

    #Add number of steps to iteration counter
    iters += m.num_steps

    if verbose and step % (epoch_size // 10) == 10:
        print("Itr %d of %d, perplexity: %.3f speed: %.
0f wps" % (step , epoch_size, np.exp(costs / iters), iters
* m.batch_size / (time.time() - start_time)))

    # Returns the Perplexity rating for us to keep track of
how the model is evolving
    return np.exp(costs / iters)

```

Now, we create the main method to tie everything together. The code here reads the data from the directory, using the reader helper module, and then trains and evaluates the model on both a testing and a validating subset of data.

In [53]:

```

# Reads the data and separates it into training data,
validation data and testing data
raw_data = reader.ptb_raw_data(data_dir)
train_data, valid_data, test_data, _, _ = raw_data

```

In [54]:

```

# Initializes the Execution Graph and the Session
with tf.Graph().as_default(), tf.Session() as session:
    initializer = tf.random_uniform_initializer(-
init_scale, init_scale)

    # Instantiates the model for training
    # tf.variable_scope add a prefix to the variables
created with tf.get_variable
    with tf.variable_scope("model", reuse=None,
initializer=initializer):
        m = PTBModel("is_training")

    # Reuses the trained parameters for the validation and
testing models
    # They are different instances but use the same

```

```

variables for weights and biases, they just don't change
when data is input
    with tf.variable_scope("model", reuse=True,
initializer=initializer):
        mvalid = PTBModel("is_validating")
        mtest = PTBModel("is_testing")

    #Initialize all variables
    tf.global_variables_initializer().run()

    for i in range(max_epoch):
        # Define the decay for this epoch
        lr_decay = decay ** max(i - max_epoch_decay_lr,
0.0)

        # Set the decayed learning rate as the learning
rate for this epoch
        m.assign_lr(session, learning_rate * lr_decay)

        print("Epoch %d : Learning rate: %.3f" % (i + 1,
session.run(m.lr)))

        # Run the loop for this epoch in the training model
        train_perplexity = run_one_epoch(session, m,
train_data, m.train_op, verbose=True)
        print("Epoch %d : Train Perplexity: %.3f" % (i + 1,
train_perplexity))

        # Run the loop for this epoch in the validation
model
        valid_perplexity = run_one_epoch(session, mvalid,
valid_data, tf.no_op())
        print("Epoch %d : Valid Perplexity: %.3f" % (i + 1,
valid_perplexity))

        # Run the loop in the testing model to see how
effective was our training
        test_perplexity = run_one_epoch(session, mtest,
test_data, tf.no_op())

        print("Test Perplexity: %.3f" % test_perplexity)

Epoch 1 : Learning rate: 1.000
Itr 10 of 774, perplexity: 4090.238 speed: 2260 wps
Itr 87 of 774, perplexity: 1273.646 speed: 2310 wps
Itr 164 of 774, perplexity: 982.987 speed: 2307 wps
Itr 241 of 774, perplexity: 819.627 speed: 2285 wps

```

Itr 318 of 774, perplexity: 724.567 speed: 2282 wps
Itr 395 of 774, perplexity: 649.213 speed: 2289 wps
Itr 472 of 774, perplexity: 587.359 speed: 2298 wps
Itr 549 of 774, perplexity: 533.067 speed: 2299 wps
Itr 626 of 774, perplexity: 489.454 speed: 2304 wps
Itr 703 of 774, perplexity: 455.080 speed: 2306 wps
Epoch 1 : Train Perplexity: 430.634
Epoch 1 : Valid Perplexity: 256.441
Epoch 2 : Learning rate: 1.000
Itr 10 of 774, perplexity: 275.627 speed: 2189 wps
Itr 87 of 774, perplexity: 238.904 speed: 2303 wps
Itr 164 of 774, perplexity: 228.938 speed: 2312 wps
Itr 241 of 774, perplexity: 219.712 speed: 2317 wps
Itr 318 of 774, perplexity: 217.133 speed: 2315 wps
Itr 395 of 774, perplexity: 211.191 speed: 2318 wps
Itr 472 of 774, perplexity: 206.899 speed: 2317 wps
Itr 549 of 774, perplexity: 200.413 speed: 2318 wps
Itr 626 of 774, perplexity: 194.814 speed: 2319 wps
Itr 703 of 774, perplexity: 190.702 speed: 2319 wps
Epoch 2 : Train Perplexity: 187.875
Epoch 2 : Valid Perplexity: 179.330
Epoch 3 : Learning rate: 1.000
Itr 10 of 774, perplexity: 188.755 speed: 2298 wps
Itr 87 of 774, perplexity: 161.439 speed: 2330 wps
Itr 164 of 774, perplexity: 157.784 speed: 2318 wps
Itr 241 of 774, perplexity: 152.883 speed: 2323 wps
Itr 318 of 774, perplexity: 152.964 speed: 2319 wps
Itr 395 of 774, perplexity: 150.248 speed: 2320 wps
Itr 472 of 774, perplexity: 148.659 speed: 2321 wps
Itr 549 of 774, perplexity: 144.914 speed: 2319 wps
Itr 626 of 774, perplexity: 142.057 speed: 2313 wps
Itr 703 of 774, perplexity: 140.223 speed: 2311 wps
Epoch 3 : Train Perplexity: 139.132
Epoch 3 : Valid Perplexity: 152.885
Epoch 4 : Learning rate: 1.000
Itr 10 of 774, perplexity: 150.212 speed: 2298 wps
Itr 87 of 774, perplexity: 128.579 speed: 2274 wps
Itr 164 of 774, perplexity: 126.476 speed: 2281 wps
Itr 241 of 774, perplexity: 123.077 speed: 2296 wps
Itr 318 of 774, perplexity: 123.639 speed: 2303 wps
Itr 395 of 774, perplexity: 121.808 speed: 2308 wps
Itr 472 of 774, perplexity: 120.991 speed: 2309 wps
Itr 549 of 774, perplexity: 118.125 speed: 2308 wps
Itr 626 of 774, perplexity: 116.220 speed: 2312 wps
Itr 703 of 774, perplexity: 115.105 speed: 2314 wps
Epoch 4 : Train Perplexity: 114.575
Epoch 4 : Valid Perplexity: 140.927

Epoch 5 : Learning rate: 1.000
Itr 10 of 774, perplexity: 127.475 speed: 2183 wps
Itr 87 of 774, perplexity: 108.904 speed: 2297 wps
Itr 164 of 774, perplexity: 107.815 speed: 2310 wps
Itr 241 of 774, perplexity: 105.222 speed: 2309 wps
Itr 318 of 774, perplexity: 105.946 speed: 2312 wps
Itr 395 of 774, perplexity: 104.478 speed: 2313 wps
Itr 472 of 774, perplexity: 104.049 speed: 2311 wps
Itr 549 of 774, perplexity: 101.734 speed: 2312 wps
Itr 626 of 774, perplexity: 100.279 speed: 2313 wps
Itr 703 of 774, perplexity: 99.505 speed: 2314 wps
Epoch 5 : Train Perplexity: 99.242
Epoch 5 : Valid Perplexity: 136.334
Epoch 6 : Learning rate: 0.500
Itr 10 of 774, perplexity: 109.463 speed: 2300 wps
Itr 87 of 774, perplexity: 93.078 speed: 2295 wps
Itr 164 of 774, perplexity: 91.362 speed: 2305 wps
Itr 241 of 774, perplexity: 88.410 speed: 2304 wps
Itr 318 of 774, perplexity: 88.570 speed: 2307 wps
Itr 395 of 774, perplexity: 86.818 speed: 2306 wps
Itr 472 of 774, perplexity: 86.060 speed: 2309 wps
Itr 549 of 774, perplexity: 83.664 speed: 2310 wps
Itr 626 of 774, perplexity: 82.051 speed: 2309 wps
Itr 703 of 774, perplexity: 81.016 speed: 2310 wps
Epoch 6 : Train Perplexity: 80.469
Epoch 6 : Valid Perplexity: 126.533
Epoch 7 : Learning rate: 0.250
Itr 10 of 774, perplexity: 95.412 speed: 2285 wps
Itr 87 of 774, perplexity: 82.098 speed: 2303 wps
Itr 164 of 774, perplexity: 80.722 speed: 2304 wps
Itr 241 of 774, perplexity: 78.096 speed: 2286 wps
Itr 318 of 774, perplexity: 78.268 speed: 2288 wps
Itr 395 of 774, perplexity: 76.613 speed: 2286 wps
Itr 472 of 774, perplexity: 75.848 speed: 2288 wps
Itr 549 of 774, perplexity: 73.613 speed: 2291 wps
Itr 626 of 774, perplexity: 72.072 speed: 2291 wps
Itr 703 of 774, perplexity: 71.011 speed: 2291 wps
Epoch 7 : Train Perplexity: 70.403
Epoch 7 : Valid Perplexity: 124.164
Epoch 8 : Learning rate: 0.125
Itr 10 of 774, perplexity: 87.555 speed: 2251 wps
Itr 87 of 774, perplexity: 75.838 speed: 2256 wps
Itr 164 of 774, perplexity: 74.743 speed: 2262 wps
Itr 241 of 774, perplexity: 72.326 speed: 2273 wps
Itr 318 of 774, perplexity: 72.535 speed: 2278 wps
Itr 395 of 774, perplexity: 70.990 speed: 2283 wps
Itr 472 of 774, perplexity: 70.259 speed: 2287 wps

Itr 549 of 774, perplexity: 68.131 speed: 2289 wps
Itr 626 of 774, perplexity: 66.643 speed: 2288 wps
Itr 703 of 774, perplexity: 65.588 speed: 2288 wps
Epoch 8 : Train Perplexity: 64.973
Epoch 8 : Valid Perplexity: 123.072
Epoch 9 : Learning rate: 0.062
Itr 10 of 774, perplexity: 83.153 speed: 2176 wps
Itr 87 of 774, perplexity: 72.360 speed: 2258 wps
Itr 164 of 774, perplexity: 71.396 speed: 2270 wps
Itr 241 of 774, perplexity: 69.130 speed: 2285 wps
Itr 318 of 774, perplexity: 69.361 speed: 2293 wps
Itr 395 of 774, perplexity: 67.884 speed: 2297 wps
Itr 472 of 774, perplexity: 67.180 speed: 2301 wps
Itr 549 of 774, perplexity: 65.124 speed: 2303 wps
Itr 626 of 774, perplexity: 63.673 speed: 2301 wps
Itr 703 of 774, perplexity: 62.632 speed: 2305 wps
Epoch 9 : Train Perplexity: 62.023
Epoch 9 : Valid Perplexity: 122.462
Epoch 10 : Learning rate: 0.031
Itr 10 of 774, perplexity: 80.891 speed: 2211 wps
Itr 87 of 774, perplexity: 70.499 speed: 2284 wps
Itr 164 of 774, perplexity: 69.555 speed: 2295 wps
Itr 241 of 774, perplexity: 67.382 speed: 2290 wps
Itr 318 of 774, perplexity: 67.617 speed: 2289 wps
Itr 395 of 774, perplexity: 66.175 speed: 2293 wps
Itr 472 of 774, perplexity: 65.481 speed: 2294 wps
Itr 549 of 774, perplexity: 63.473 speed: 2296 wps
Itr 626 of 774, perplexity: 62.043 speed: 2294 wps
Itr 703 of 774, perplexity: 61.019 speed: 2297 wps
Epoch 10 : Train Perplexity: 60.417
Epoch 10 : Valid Perplexity: 122.073
Epoch 11 : Learning rate: 0.016
Itr 10 of 774, perplexity: 79.714 speed: 2288 wps
Itr 87 of 774, perplexity: 69.459 speed: 2276 wps
Itr 164 of 774, perplexity: 68.535 speed: 2274 wps
Itr 241 of 774, perplexity: 66.402 speed: 2274 wps
Itr 318 of 774, perplexity: 66.637 speed: 2289 wps
Itr 395 of 774, perplexity: 65.223 speed: 2295 wps
Itr 472 of 774, perplexity: 64.534 speed: 2299 wps
Itr 549 of 774, perplexity: 62.555 speed: 2302 wps
Itr 626 of 774, perplexity: 61.140 speed: 2299 wps
Itr 703 of 774, perplexity: 60.128 speed: 2299 wps
Epoch 11 : Train Perplexity: 59.533
Epoch 11 : Valid Perplexity: 121.913
Epoch 12 : Learning rate: 0.008
Itr 10 of 774, perplexity: 79.105 speed: 2242 wps
Itr 87 of 774, perplexity: 68.870 speed: 2315 wps

Itr 164 of 774, perplexity: 67.972 speed: 2315 wps
Itr 241 of 774, perplexity: 65.861 speed: 2320 wps
Itr 318 of 774, perplexity: 66.092 speed: 2316 wps
Itr 395 of 774, perplexity: 64.697 speed: 2310 wps
Itr 472 of 774, perplexity: 64.013 speed: 2313 wps
Itr 549 of 774, perplexity: 62.051 speed: 2311 wps
Itr 626 of 774, perplexity: 60.643 speed: 2311 wps
Itr 703 of 774, perplexity: 59.637 speed: 2314 wps
Epoch 12 : Train Perplexity: 59.045
Epoch 12 : Valid Perplexity: 121.826
Epoch 13 : Learning rate: 0.004
Itr 10 of 774, perplexity: 78.787 speed: 2279 wps
Itr 87 of 774, perplexity: 68.557 speed: 2306 wps
Itr 164 of 774, perplexity: 67.671 speed: 2302 wps
Itr 241 of 774, perplexity: 65.576 speed: 2308 wps
Itr 318 of 774, perplexity: 65.805 speed: 2303 wps
Itr 395 of 774, perplexity: 64.415 speed: 2309 wps
Itr 472 of 774, perplexity: 63.735 speed: 2306 wps
Itr 549 of 774, perplexity: 61.783 speed: 2298 wps
Itr 626 of 774, perplexity: 60.378 speed: 2299 wps
Itr 703 of 774, perplexity: 59.374 speed: 2294 wps
Epoch 13 : Train Perplexity: 58.785
Epoch 13 : Valid Perplexity: 121.706
Epoch 14 : Learning rate: 0.002
Itr 10 of 774, perplexity: 78.599 speed: 2224 wps
Itr 87 of 774, perplexity: 68.389 speed: 2288 wps
Itr 164 of 774, perplexity: 67.510 speed: 2296 wps
Itr 241 of 774, perplexity: 65.426 speed: 2297 wps
Itr 318 of 774, perplexity: 65.657 speed: 2296 wps
Itr 395 of 774, perplexity: 64.269 speed: 2297 wps
Itr 472 of 774, perplexity: 63.590 speed: 2299 wps
Itr 549 of 774, perplexity: 61.643 speed: 2301 wps
Itr 626 of 774, perplexity: 60.240 speed: 2298 wps
Itr 703 of 774, perplexity: 59.238 speed: 2299 wps
Epoch 14 : Train Perplexity: 58.649
Epoch 14 : Valid Perplexity: 121.569
Epoch 15 : Learning rate: 0.001
Itr 10 of 774, perplexity: 78.476 speed: 2286 wps
Itr 87 of 774, perplexity: 68.285 speed: 2290 wps
Itr 164 of 774, perplexity: 67.414 speed: 2300 wps
Itr 241 of 774, perplexity: 65.338 speed: 2300 wps
Itr 318 of 774, perplexity: 65.573 speed: 2300 wps
Itr 395 of 774, perplexity: 64.188 speed: 2299 wps
Itr 472 of 774, perplexity: 63.511 speed: 2301 wps
Itr 549 of 774, perplexity: 61.567 speed: 2303 wps
Itr 626 of 774, perplexity: 60.165 speed: 2301 wps
Itr 703 of 774, perplexity: 59.165 speed: 2296 wps

Epoch 15 : Train Perplexity: 58.577
Epoch 15 : Valid Perplexity: 121.474
Test Perplexity: 117.999

As you can see, the model's perplexity rating drops very quickly after a few iterations. As was elaborated before, **lower Perplexity means that the model is more certain about its prediction**. As such, we can be sure that this model is performing well!

This is the end of the **Applying Recurrent Neural Networks to Text Processing** notebook. Hopefully you now have a better understanding of Recurrent Neural Networks and how to implement one utilizing TensorFlow. Thank you for reading this notebook, and good luck on your studies.

Want to learn more?¶

Running deep learning programs usually needs a high performance platform. **PowerAI** speeds up deep learning and AI. Built on IBM's Power Systems, **PowerAI** is a scalable software platform that accelerates deep learning and AI with blazing performance for individual users or enterprises. The **PowerAI** platform supports popular machine learning libraries and dependencies including TensorFlow, Caffe, Torch, and Theano. You can use [PowerAI on IMB Cloud](#).

Also, you can use **Watson Studio** to run these notebooks faster with bigger datasets. **Watson Studio** is IBM's leading cloud solution for data scientists, built by data scientists. With Jupyter notebooks, RStudio, Apache Spark and popular libraries pre-packaged in the cloud, **Watson Studio** enables data scientists to collaborate on their projects without having to install anything. Join the fast-growing community of **Watson Studio** users today with a free account at [Watson Studio](#). This is the end of this lesson. Thank you for reading this notebook, and good luck on your studies.

Thanks for completing this lesson!¶

Notebook created by [Walter Gomes de Amorim Junior](#), <a href = "<https://linkedin.com/in/saeedaghabozorgi>"> Saeed Aghabozorgi </h4>

Copyright © 2018 [Cognitive Class](#). This notebook and its source code are released under the terms of the [MIT License](#).

In []:

