WILDML

Artificial Intelligence, Deep Learning, and NLP

SEPTEMBER 30, 2015 BY DENNY BRITZ

This the second part of the Recurrent Neural Network Tutorial. The first part is here.

Code to follow along is on Github.

In this part we will implement a full Recurrent Neural Network from scratch using Python and optimize our implementation using Theano, a library to perform operations on a GPU. **The full code is available on Github**. I will skip over some boilerplate code that is not essential to understanding Recurrent Neural Networks, but all of that is also on Github.

Our goal is to build a Language Model using a Recurrent Neural Network. Here's what that means. Let's say we have sentence of m words. A language model allows us to predict the probability of observing the sentence (in a given dataset) as:

$$P(w_1,...,w_m) = \prod_{i=1}^m P(w_i \mid w_1,...,w_{i-1})$$

In words, the probability of a sentence is the product of probabilities of each word given the words that came before it. So, the probability of the sentence "He went to buy some chocolate" would be the probability of "chocolate" given "He went to buy some", multiplied by the probability of "some" given "He went to buy", and so on.

Why is that useful? Why would we want to assign a probability to observing a sentence?

First, such a model can be used as a scoring mechanism. For example, a Machine Translation system typically generates multiple candidates for an input sentence. You could use a language model to pick the most probable sentence. Intuitively, the most probable sentence is likely to be grammatically correct. Similar scoring happens in speech recognition systems.

But solving the Language Modeling problem also has a cool side effect. Because we can predict the probability of a word given the preceding words, we are able to generate new text. It's a . Given an existing sequence of words we sample a next word from the predicted probabilities, and repeat the process until we have a full sentence. Andrej Karparthy has a great post that demonstrates what language models are capable of. His models are trained on single characters as opposed to full words, and can generate anything from Shakespeare to Linux Code.

Note that in the above equation the probability of each word is conditioned on **all** previous words. In practice, many models have a hard time representing such long-term dependencies due to computational or memory constraints. They are typically limited to looking at only a few of the previous words. RNNs can, in theory, capture such long-term dependencies, but in practice it's a bit more complex. We'll explore that in a later post.

To train our language model we need text to learn from. Fortunately we don't need any labels to train a language model, just raw text. I downloaded 15,000 longish reddit comments from a dataset available on Google's BigQuery. Text generated by our model will sound like reddit commenters (hopefully)! But as with most Machine Learning projects we first need to do some pre-processing to get our data into the right format.

1. Tokenize Text

We have raw text, but we want to make predictions on a per-word basis. This means we must our comments into sentences, and sentences into words. We could just split each of the comments by spaces, but that wouldn't handle punctuation properly. The sentence "He left!" should be 3 tokens: "He", "left", "!". We'll use NLTK's word_tokenize and sent tokenize methods, which do most of the hard work for us.

2. Remove infrequent words

Most words in our text will only appear one or two times. It's a good idea to remove these infrequent words. Having a huge vocabulary will make our model slow to train (we'll talk about why that is later), and because we don't have a lot of contextual examples for such words we wouldn't be able to learn how to use them correctly anyway. That's quite similar to how humans learn. To really understand how to appropriately use a word you need to have seen it in different contexts.

In our code we limit our vocabulary to the vocabulary_size most common words (which I set to 8000, but feel free to change it). We replace all words not included in our vocabulary by UNKNOWN_TOKEN. For example, if we don't include the word "nonlinearities" in our vocabulary, the sentence "nonlineraties are important in neural networks" becomes "UNKNOWN_TOKEN are important in Neural Networks". The word UNKNOWN_TOKEN will become part of our vocabulary and we will predict it just like any other word. When we generate new text we can replace UNKNOWN_TOKEN again, for example by taking a randomly sampled word not in our vocabulary, or we could just generate sentences until we get one that doesn't contain an unknown token.

3. Prepend special start and end tokens

We also want to learn which words tend start and end a sentence. To do this we prepend a special SENTENCE_START token, and append a special SENTENCE_END token to each sentence. This allows us to ask: Given that the first token is SENTENCE_START, what is the likely next word (the actual first word of the sentence)?

4. Build training data matrices

The input to our Recurrent Neural Networks are vectors, not strings. So we create a mapping between words and indices, index_to_word, and word_to_index. For example, the word "friendly" may be at index 2001. A training example *x* may look like [0, 179, 341,

416], where O corresponds to SENTENCE_START. The corresponding label y would be [179, 341, 416, 1]. Remember that our goal is to predict the next word, so y is just the x vector shifted by one position with the last element being the SENTENCE_END token. In other words, the correct prediction for word 179 above would be 341, the actual next word.

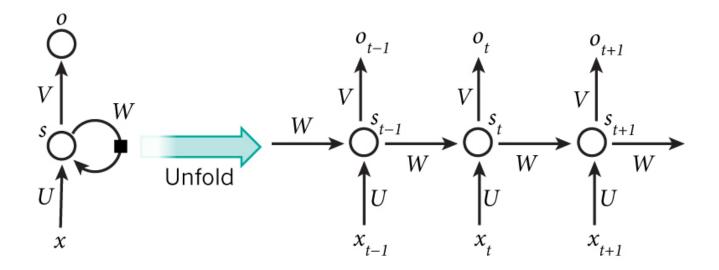
```
1 vocabulary_size = 8000
   unknown_token = "UNKNOWN_TOKEN"
  sentence_start_token = "SENTENCE_START"
   sentence_end_token = "SENTENCE_END"
 5
 6 # Read the data and append SENTENCE_START and SENTENCE_END tokens
 7
   print "Reading CSV file..."
   with open('data/reddit-comments-2015-08.csv', 'rb') as f:
 9
       reader = csv.reader(f, skipinitialspace=True)
10
       reader.next()
       # Split full comments into sentences
11
12
       sentences = itertools.chain(*[nltk.sent_tokenize(x[0].decode('utf-8').lower()) for
13 x in reader])
14
       # Append SENTENCE_START and SENTENCE_END
15
       sentences = ["%s %s %s" % (sentence_start_token, x, sentence_end_token) f
16 or x in sentences
17 | print " Parsed %d sentences. " % (len(sentences))
18
19
   # Tokenize the sentences into words
20 tokenized_sentences = [nltk.word_tokenize(sent) for sent in sentences]
21
22
   # Count the word frequencies
   word_freq = nltk.FreqDist(itertools.chain(*tokenized_sentences))
   print "Found %d unique words tokens." % len(word_freq.items())
25
26 | # Get the most common words and build index_to_word and word_to_index vectors
27 | vocab = word_freq.most_common(vocabulary_size-1)
28 index_to_word = [x[0] for x in vocab]
29 index_to_word.append(unknown_token)
30 word_to_index = dict([(w,i) for i,w in enumerate(index_to_word)])
31
32 print "Using vocabulary size %d." % vocabulary_size
33 print "The least frequent word in our vocabulary is '%s' and appeared %d times.&qu
34 ot; % (vocab[-1][0], vocab[-1][1])
35
36 # Replace all words not in our vocabulary with the unknown token
37 for i, sent in enumerate(tokenized_sentences):
38
       tokenized_sentences[i] = [w if w in word_to_index else unknown_token for w in sent]
40 print "\nExample sentence: '%s'" % sentences[0]
41 print " \nExample sentence after Pre-processing: '%s'" % tokenized_sentences[0
42 7
   # Create the training data
   X_train = np.asarray([[word_to_index[w] for w in sent[:-1]] for sent in tokenized_sente
   y_train = np.asarray([[word_to_index[w] for w in sent[1:]] for sent in tokenized_senten
   ces])
```

Here's an actual training example from our text:

```
x:
SENTENCE_START what are n't you understanding about this ? !
[0, 51, 27, 16, 10, 856, 53, 25, 34, 69]

y:
what are n't you understanding about this ? ! SENTENCE_END
[51, 27, 16, 10, 856, 53, 25, 34, 69, 1]
```

For a general overview of RNNs take a look at first part of the tutorial.



Let's get concrete and see what the RNN for our language model looks like. The input x will be a sequence of words (just like the example printed above) and each x_t is a single word. But there's one more thing: Because of how matrix multiplication works we can't simply use a word index (like 36) as an input. Instead, we represent each word as a

of size vocabulary_size. For example, the word with index 36 would be the vector of all 0's and a 1 at position 36. So, each x_t will become a vector, and x will be a matrix, with each row representing a word. We'll perform this transformation in our Neural Network code instead of doing it in the pre-processing. The output of our network o has a similar format. Each o_t is a vector of vocabulary_size elements, and each element represents the probability of that word being the next word in the sentence.

Let's recap the equations for the RNN from the first part of the tutorial:

$$s_t = \tanh(Ux_t + Ws_{t-1})$$

 $o_t = \operatorname{softmax}(Vs_t)$

I always find it useful to write down the dimensions of the matrices and vectors. Let's assume we pick a vocabulary size C=8000 and a hidden layer size H=100. You can think of the hidden layer size as the "memory" of our network. Making it bigger allows us to learn more complex patterns, but also results in additional computation. Then we have:

$$x_t \in \mathbb{R}^{8000}$$

 $o_t \in \mathbb{R}^{8000}$
 $s_t \in \mathbb{R}^{100}$
 $U \in \mathbb{R}^{100 \times 8000}$
 $V \in \mathbb{R}^{8000 \times 100}$
 $W \in \mathbb{R}^{100 \times 100}$

This is valuable information. Remember that U,V and W are the parameters of our network we want to learn from data. Thus, we need to learn a total of $2HC+H^2$ parameters. In the case of C=8000 and H=100 that's 1,610,000. The dimensions also tell us the bottleneck of our model. Note that because x_t is a one-hot vector, multiplying it with W is essentially the same as selecting a column of W, so we don't need to perform the full multiplication. Then, the biggest matrix multiplication in our network is W that's why we want to keep our vocabulary size small if possible.

Armed with this, it's time to start our implementation.

Initialization

We start by declaring a RNN class an initializing our parameters. I'm calling this class RNNNumpy because we will implement a Theano version later. Initializing the parameters U,V and W is a bit tricky. We can't just initialize them to 0's because that would result in symmetric calculations in all our layers. We must initialize them randomly. Because proper initialization seems to have an impact on training results there has been lot of research in this area. It turns out that the best initialization depends on the activation function (\tanh in our case) and one recommended approach is to initialize the weights randomly in the interval from $\begin{bmatrix} -\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \end{bmatrix}$ where n is the number of incoming connections from the previous

layer. This may sound overly complicated, but don't worry too much it. As long as you initialize your parameters to small random values it typically works out fine.

```
class RNNNumpy:
1
2
3
       def __init__(self, word_dim, hidden_dim=100, bptt_truncate=4):
4
           # Assign instance variables
5
           self.word_dim = word_dim
6
           self.hidden_dim = hidden_dim
7
           self.bptt_truncate = bptt_truncate
8
           # Randomly initialize the network parameters
           self.U = np.random.uniform(-np.sqrt(1./word_dim), np.sqrt(1./word_dim), (hidden
9
10
   _dim, word_dim))
           self.V = np.random.uniform(-np.sqrt(1./hidden_dim), np.sqrt(1./hidden_dim), (wo
   rd_dim, hidden_dim))
           self.W = np.random.uniform(-np.sqrt(1./hidden_dim), np.sqrt(1./hidden_dim), (hi
   dden_dim, hidden_dim))
```

Above, word_dim is the size of our vocabulary, and hidden_dim is the size of our hidden layer (we can pick it). Don't worry about the bptt_truncate parameter for now, we'll explain what that is later.

Forward Propagation

Next, let's implement the forward propagation (predicting word probabilities) defined by our equations above:

```
def forward_propagation(self, x):
2
       # The total number of time steps
3
       T = len(x)
4
       # During forward propagation we save all hidden states in s because need them late
5
   r.
6
       # We add one additional element for the initial hidden, which we set to 0
7
       s = np.zeros((T + 1, self.hidden_dim))
8
       s[-1] = np.zeros(self.hidden_dim)
9
       # The outputs at each time step. Again, we save them for later.
10
       o = np.zeros((T, self.word_dim))
11
       # For each time step...
12
       for t in np.arange(T):
13
           # Note that we are indxing U by x[t]. This is the same as multiplying U with a
14
    one-hot vector.
15
           s[t] = np.tanh(self.U[:,x[t]] + self.W.dot(s[t-1]))
16
           o[t] = softmax(self.V.dot(s[t]))
17
       return [o, s]
   RNNNumpy.forward_propagation = forward_propagation
```

We not only return the calculated outputs, but also the hidden states. We will use them later to calculate the gradients, and by returning them here we avoid duplicate computation. Each o_t is a vector of probabilities representing the words in our vocabulary,

but sometimes, for example when evaluating our model, all we want is the next word with the highest probability. We call this function predict:

```
def predict(self, x):
    # Perform forward propagation and return index of the highest score
    o, s = self.forward_propagation(x)
    return np.argmax(o, axis=1)

RNNNumpy.predict = predict
```

Let's try our newly implemented methods and see an example output:

```
np.random.seed(10)
model = RNNNumpy(vocabulary_size)
o, s = model.forward_propagation(X_train[10])
print o.shape
print o
```

```
(45, 8000)
[[ 0.00012408  0.0001244
                          0.00012603 ...,
                                           0.00012515
                                                       0.00012488
  0.00012508]
[ 0.00012536
              0.00012582 0.00012436 ...,
                                                       0.00012456
                                           0.00012482
  0.000124511
[ 0.00012387  0.0001252
                          0.00012474 ...,
                                           0.00012559 0.00012588
  0.00012551]
 [ 0.00012414  0.00012455
                         0.0001252
                                           0.00012487
                                                       0.00012494
  0.0001263 ]
[ 0.0001252
              0.00012393 0.00012509 ...,
                                           0.00012407
                                                       0.00012578
  0.00012502]
[ 0.00012472  0.0001253
                          0.00012487 ..., 0.00012463 0.00012536
  0.0001266511
```

For each word in the sentence (45 above), our model made 8000 predictions representing probabilities of the next word. Note that because we initialized U, V, W to random values these predictions are completely random right now. The following gives the indices of the highest probability predictions for each word:

```
1 predictions = model.predict(X_train[10])
2 print predictions.shape
3 print predictions
```

```
(45,)
[1284 5221 7653 7430 1013 3562 7366 4860 2212 6601 7299 4556 2481 238 2!
21 6548 261 1780 2005 1810 5376 4146 477 7051 4832 4991 897 3485 21
7291 2007 6006 760 4864 2182 6569 2800 2752 6821 4437 7021 7875 6912 3!
```

Calculating the Loss

To train our network we need a way to measure the errors it makes. We call this the loss function L, and our goal is find the parameters U, V and W that minimize the loss function for our training data. A common choice for the loss function is the cross-entropy loss. If we have N training examples (words in our text) and C classes (the size of our vocabulary) then the loss with respect to our predictions o and the true labels y is given by:

$$L(y, o) = -\frac{1}{N} \sum_{n \in N} y_n \log o_n$$

The formula looks a bit complicated, but all it really does is sum over our training examples and add to the loss based on how off our prediction are. The further away y (the correct words) and o (our predictions), the greater the loss will be. We implement the function calculate_loss:

```
def calculate_total_loss(self, x, y):
 1
 2
       L = 0
 3
       # For each sentence...
 4
       for i in np.arange(len(y)):
 5
           o, s = self.forward_propagation(x[i])
 6
           # We only care about our prediction of the "correct" words
 7
           correct_word_predictions = o[np.arange(len(y[i])), y[i]]
 8
           # Add to the loss based on how off we were
 9
           L += -1 * np.sum(np.log(correct_word_predictions))
10
       return L
11
12
   def calculate_loss(self, x, y):
       # Divide the total loss by the number of training examples
13
14
       N = np.sum((len(y_i) for y_i in y))
15
       return self.calculate_total_loss(x,y)/N
16
17 RNNNumpy.calculate_total_loss = calculate_total_loss
18 RNNNumpy.calculate_loss = calculate_loss
```

Let's take a step back and think about what the loss should be for random predictions. That will give us a baseline and make sure our implementation is correct. We have C

words in our vocabulary, so each word should be (on average) predicted with probability 1/C, which would yield a loss of $L = -\frac{1}{N}N\log\frac{1}{C} = \log C$:

```
1 # Limit to 1000 examples to save time
2 print "Expected Loss for random predictions: %f" % np.log(vocabulary_size)
3 print "Actual loss: %f" % model.calculate_loss(X_train[:1000], y_train[:1000])
```

```
Expected Loss for random predictions: 8.987197 Actual loss: 8.987440
```

Pretty close! Keep in mind that evaluating the loss on the full dataset is an expensive operation and can take hours if you have a lot of data!

Training the RNN with SGD and Backpropagation Through Time (BPTT)

Remember that we want to find the parameters U, V and W that minimize the total loss on the training data. The most common way to do this is SGD, Stochastic Gradient Descent. The idea behind SGD is pretty simple. We iterate over all our training examples and during each iteration we nudge the parameters into a direction that reduces the error. These directions are given by the gradients on the loss: $\frac{\partial L}{\partial U}, \frac{\partial L}{\partial V}, \frac{\partial L}{\partial W}$. SGD also needs a , which defines how big of a step we want to make in each iteration. SGD is the most popular optimization method not only for Neural Networks, but also for many other Machine Learning algorithms. As such there has been a lot of research on how to optimize SGD using batching, parallelism and adaptive learning rates. Even though the basic idea is simple, implementing SGD in a really efficient way can become very complex. If you want to learn more about SGD $\underline{\text{this}}$ is a good place to start. Due to its popularity there are a wealth of tutorials floating around the web, and I don't want to duplicate them here. I'll implement a simple version of SGD that should be understandable even without a background in optimization.

But how do we calculate those gradients we mentioned above? In a traditional Neural Network we do this through the backpropagation algorithm. In RNNs we use a slightly modified version of the this algorithm called .

Because the parameters are shared by all time steps in the network, the gradient at each output depends not only on the calculations of the current time step, but also the previous time steps. If you know calculus, it really is just applying the chain rule. The next part of the tutorial will be all about BPTT, so I won't go into detailed derivation here. For a general

introduction to backpropagation check out this and this post. For now you can treat BPTT as a black box. It takes as input a training example (x, y) and returns the gradients $\frac{\partial L}{\partial U}, \frac{\partial L}{\partial V}, \frac{\partial L}{\partial W}$.

```
def bptt(self, x, y):
 2
       T = len(y)
       # Perform forward propagation
 4
       o, s = self.forward_propagation(x)
 5
       # We accumulate the gradients in these variables
 6
       dLdU = np.zeros(self.U.shape)
 7
       dLdV = np.zeros(self.V.shape)
 8
       dLdW = np.zeros(self.W.shape)
 9
       delta_o = o
       delta_o[np.arange(len(y)), y] = 1.
10
       # For each output backwards...
11
12
       for t in np.arange(T)[::-1]:
            dLdV += np.outer(delta_o[t], s[t].T)
13
            # Initial delta calculation
14
15
            delta_t = self.V.T.dot(delta_o[t]) * (1 - (s[t] ** 2))
            # Backpropagation through time (for at most self.bptt_truncate steps)
16
17
            for bptt_step in np.arange(max(0, t-self.bptt_truncate), t+1)[::-1]:
18
                # print "Backpropagation step t=%d bptt step=%d "                  % (t, bptt_ste
19 p)
20
                dLdW += np.outer(delta_t, s[bptt_step-1])
                dLdU[:,x[bptt_step]] += delta_t
21
22
                # Update delta for next step
23
                delta_t = self.W.T.dot(delta_t) * (1 - s[bptt_step-1] ** 2)
24
       return [dLdU, dLdV, dLdW]
25
   RNNNumpy.bptt = bptt
```

Gradient Checking

Whenever you implement backpropagation it is good idea to also implement

, which is a way of verifying that your implementation is correct. The idea behind gradient checking is that derivative of a parameter is equal to the slope at the point, which we can approximate by slightly changing the parameter and then dividing by the change:

$$\frac{\partial L}{\partial \theta} \approx \lim_{h \to 0} \frac{J(\theta + h) - J(\theta - h)}{2h}$$

We then compare the gradient we calculated using backpropagation to the gradient we estimated with the method above. If there's no large difference we are good. The approximation needs to calculate the total loss for parameter, so that gradient checking is very expensive (remember, we had more than a million parameters in the example above). So it's a good idea to perform it on a model with a smaller vocabulary.

```
1 def gradient_check(self, x, y, h=0.001, error_threshold=0.01):
```

```
# Calculate the gradients using backpropagation. We want to checker if these are co
4
   rrect.
5
       bptt_gradients = self.bptt(x, y)
6
       # List of all parameters we want to check.
7
       model_parameters = ['U', 'V', 'W']
8
       # Gradient check for each parameter
9
       for pidx, pname in enumerate(model_parameters):
10
           # Get the actual parameter value from the mode, e.g. model.W
11
           parameter = operator.attrgetter(pname)(self)
12
           print " Performing gradient check for parameter %s with size %d." % (p
13
   name, np.prod(parameter.shape))
14
           # Iterate over each element of the parameter matrix, e.g. (0,0), (0,1), ...
15
           it = np.nditer(parameter, flags=['multi_index'], op_flags=['readwrite'])
           while not it.finished:
16
17
               ix = it.multi_index
18
               # Save the original value so we can reset it later
19
               original_value = parameter[ix]
20
               # Estimate the gradient using (f(x+h) - f(x-h))/(2*h)
21
               parameter[ix] = original_value + h
22
               gradplus = self.calculate_total_loss([x],[y])
23
               parameter[ix] = original_value - h
24
               gradminus = self.calculate_total_loss([x],[y])
25
               estimated_gradient = (gradplus - gradminus)/(2*h)
26
               # Reset parameter to original value
27
               parameter[ix] = original_value
28
               # The gradient for this parameter calculated using backpropagation
29
               backprop_gradient = bptt_gradients[pidx][ix]
30
               # calculate The relative error: (|x - y|/(|x| + |y|))
31
               relative_error = np.abs(backprop_gradient - estimated_gradient)/(np.abs(bac
32
   kprop_gradient) + np.abs(estimated_gradient))
33
               # If the error is to large fail the gradient check
               if relative_error > error_threshold:
34
35
                   print " Gradient Check ERROR: parameter=%s ix=%s" % (pname, ix
36
   )
37
                   print "+h Loss: %f" % gradplus
38
                   print "-h Loss: %f" % gradminus
39
                   print "Estimated_gradient: %f" % estimated_gradient
40
                   print " Backpropagation gradient: %f" % backprop_gradient
41
                   print "Relative Error: %f" % relative_error
42
                   return
43
               it.iternext()
           print " Gradient check for parameter %s passed. " % (pname)
44
45
46 RNNNumpy.gradient_check = gradient_check
47
   # To avoid performing millions of expensive calculations we use a smaller vocabulary si
   ze for checking.
   grad_check_vocab_size = 100
   np.random.seed(10)
   model = RNNNumpy(grad_check_vocab_size, 10, bptt_truncate=1000)
   model.gradient_check([0,1,2,3], [1,2,3,4])
```

SGD Implementation

Now that we are able to calculate the gradients for our parameters we can implement SGD. I like to do this in two steps: 1. A function sdg_step that calculates the gradients and performs the updates for one batch. 2. An outer loop that iterates through the training set and adjusts the learning rate.

```
1 # Performs one step of SGD.
   def numpy_sdg_step(self, x, y, learning_rate):
       # Calculate the gradients
3
4
       dLdU, dLdV, dLdW = self.bptt(x, y)
5
       # Change parameters according to gradients and learning rate
6
       self.U -= learning_rate * dLdU
7
       self.V -= learning_rate * dLdV
8
       self.W -= learning_rate * dLdW
9
10 RNNNumpy.sgd_step = numpy_sdg_step
```

```
1 # Outer SGD Loop
2
   # - model: The RNN model instance
3 # - X_train: The training data set
4 # - y_train: The training data labels
5 # - learning_rate: Initial learning rate for SGD
  # - nepoch: Number of times to iterate through the complete dataset
7
   # - evaluate_loss_after: Evaluate the loss after this many epochs
   def train_with_sqd(model, X_train, y_train, learning_rate=0.005, nepoch=100, evaluate_l
9
   oss_after=5):
10
       # We keep track of the losses so we can plot them later
11
       losses = □
12
       num_examples_seen = 0
13
       for epoch in range(nepoch):
14
           # Optionally evaluate the loss
15
           if (epoch % evaluate_loss_after == 0):
               loss = model.calculate_loss(X_train, y_train)
16
17
               losses.append((num_examples_seen, loss))
               time = datetime.now().strftime('%Y-%m-%d %H:%M:%S')
18
19
               print "%s: Loss after num_examples_seen=%d epoch=%d: %f" % (time,
20 num_examples_seen, epoch, loss)
21
               # Adjust the learning rate if loss increases
22
               if (len(losses) > 1 and losses[-1][1] > losses[-2][1]):
23
                   learning_rate = learning_rate * 0.5
24
                   print " Setting learning rate to %f" % learning_rate
25
               sys.stdout.flush()
26
           # For each training example...
27
           for i in range(len(y_train)):
28
               # One SGD step
               model.sqd_step(X_train[i], y_train[i], learning_rate)
               num_examples_seen += 1
```

Done! Let's try to get a sense of how long it would take to train our network:

```
1  np.random.seed(10)
2  model = RNNNumpy(vocabulary_size)
3  %timeit model.sgd_step(X_train[10], y_train[10], 0.005)
```

Uh-oh, bad news. One step of SGD takes approximately 350 milliseconds on my laptop. We have about 80,000 examples in our training data, so one epoch (iteration over the whole data set) would take several hours. Multiple epochs would take days, or even weeks! And we're still working with a small dataset compared to what's being used by many of the companies and researchers out there. What now?

Fortunately there are many ways to speed up our code. We could stick with the same model and make our code run faster, or we could modify our model to be less computationally expensive, or both. Researchers have identified many ways to make models less computationally expensive, for example by using a hierarchical softmax or adding projection layers to avoid the large matrix multiplications (see also here or here). But I want to keep our model simple and go the first route: Make our implementation run faster using a GPU. Before doing that though, let's just try to run SGD with a small dataset and check if the loss actually decreases:

```
1    np.random.seed(10)
2  # Train on a small subset of the data to see what happens
3  model = RNNNumpy(vocabulary_size)
4  losses = train_with_sgd(model, X_train[:100], y_train[:100], nepoch=10, evaluate_loss_af ter=1)
```

```
2015-09-30 10:08:19: Loss after num_examples_seen=0 epoch=0: 8.987425 2015-09-30 10:08:35: Loss after num_examples_seen=100 epoch=1: 8.976270 2015-09-30 10:08:50: Loss after num_examples_seen=200 epoch=2: 8.960212 2015-09-30 10:09:06: Loss after num_examples_seen=300 epoch=3: 8.930430 2015-09-30 10:09:22: Loss after num_examples_seen=400 epoch=4: 8.862264 2015-09-30 10:09:38: Loss after num_examples_seen=500 epoch=5: 6.913570 2015-09-30 10:09:53: Loss after num_examples_seen=600 epoch=6: 6.302493 2015-09-30 10:10:07: Loss after num_examples_seen=700 epoch=7: 6.014995 2015-09-30 10:10:24: Loss after num_examples_seen=800 epoch=8: 5.833877 2015-09-30 10:10:39: Loss after num_examples_seen=900 epoch=9: 5.710718
```

Good, it seems like our implementation is at least doing something useful and decreasing the loss, just like we wanted.

I have previously written a tutorial on Theano, and since all our logic will stay exactly the same I won't go through optimized code here again. I defined a RNNTheano class that replaces the numpy calculations with corresponding calculations in Theano. Just like the rest of this post, the code is also available Github.

```
1 | np.random.seed(10)
```

```
2 model = RNNTheano(vocabulary_size)
3 %timeit model.sgd_step(X_train[10], y_train[10], 0.005)
```

This time, one SGD step takes 70ms on my Mac (without GPU) and 23ms on a g2.2xlarge Amazon EC2 instance with GPU. That's a 15x improvement over our initial implementation and means we can train our model in hours/days instead of weeks. There are still a vast number of optimizations we could make, but we're good enough for now.

To help you avoid spending days training a model I have pre-trained a Theano model with a hidden layer dimensionality of 50 and a vocabulary size of 8000. I trained it for 50 epochs in about 20 hours. The loss was was still decreasing and training longer would probably have resulted in a better model, but I was running out of time and wanted to publish this post. Feel free to try it out yourself and trian for longer. You can find the model parameters in data/trained-model-theano.npz in the Github repository and load them using the load_model_parameters_theano method:

```
from utils import load_model_parameters_theano, save_model_parameters_theano
model = RNNTheano(vocabulary_size, hidden_dim=50)
# losses = train_with_sgd(model, X_train, y_train, nepoch=50)
# save_model_parameters_theano('./data/trained-model-theano.npz', model)
load_model_parameters_theano('./data/trained-model-theano.npz', model)
```

Now that we have our model we can ask it to generate new text for us! Let's implement a helper function to generate new sentences:

```
def generate_sentence(model):
2
       # We start the sentence with the start token
3
       new_sentence = [word_to_index[sentence_start_token]]
4
       # Repeat until we get an end token
5
6
       while not new_sentence[-1] == word_to_index[sentence_end_token]:
           next_word_probs = model.forward_propagation(new_sentence)
7
           sampled_word = word_to_index[unknown_token]
8
           # We don't want to sample unknown words
9
           while sampled_word == word_to_index[unknown_token]:
10
                samples = np.random.multinomial(1, next_word_probs[-1])
11
                sampled_word = np.argmax(samples)
12
           new_sentence.append(sampled_word)
13
       sentence_str = [index_to_word[x] for x in new_sentence[1:-1]]
14
       return sentence_str
15
16 num_sentences = 10
17
   senten_min_length = 7
18
19 for i in range(num_sentences):
```

```
20    sent = []
21    # We want long sentences, not sentences with one or two words
22    while len(sent) & amp; lt; senten_min_length:
23         sent = generate_sentence(model)
24    print & quot; & quot; .join(sent)
```

A few selected (censored) sentences. I added capitalization.

- Anyway, to the city scene you're an idiot teenager.
- What ?!!!!ignore!
- Screw fitness, you're saying: https
- Thanks for the advice to keep my thoughts around girls.
- Yep, please disappear with the terrible generation.

Looking at the generated sentences there are a few interesting things to note. The model successfully learn syntax. It properly places commas (usually before and's and or's) and ends sentence with punctuation. Sometimes it mimics internet speech such as multiple exclamation marks or smileys.

However, the vast majority of generated sentences don't make sense or have grammatical errors (I really picked the best ones above). One reason could be that we did not train our network long enough (or didn't use enough training data). That may be true, but it's most likely not the main reason. **Our vanilla RNN can't generate meaningful text because it's unable to learn dependencies between words that are several steps apart**. That's also why RNNs failed to gain popularity when they were first invented. They were beautiful in theory but didn't work well in practice, and we didn't immediately understand why.

Fortunately, the difficulties in training RNNs are much better understood now. In the next part of this tutorial we will explore the Backpropagation Through Time (BPTT) algorithm in more detail and demonstrate what's called the . This will motivate our move to more sophisticated RNN models, such as LSTMs, which are the current state of the art for many tasks in NLP (and can generate much better reddit comments!). Everything you learned in this tutorial also applies to LSTMs and other RNN models, so don't feel discouraged if the results for a vanilla RNN are worse then you expected.

That's it for now. **Please leave questions or feedback in the comments!** and don't forget to check out the [code 1="Github" language="on"]/code.

DEEP LEARNING, GPU, LANGUAGE MODELING, RECURRENT NEURAL NETWORKS





tiendh • 3 years ago

Helo Denny,

Thanks for your great post. I have a question. To save training time, I'd like to use your trained model. I load the model by the following code:

model = RNNTheano(vocabulary_size, hidden_dim=_HIDDEN_DIM) load_model_parameters_theano(_MODEL_FILE, model)

I want to predict the next word of a test sentence like: 'hello, how are' which I expect the next word is 'you', so from the test sentence and the same vocabulary, I create the input vector:

input feature vector = [6538 5631 2366 3148].

Next, I use function 'predict' of the model as follows: predicted_feature_vector = model.predict(input_feature_vector)

And it throws the error:

Traceback (most recent call last):

File "test_trained_model.py", line 67, in <module>

predicted feature vector = model.predict(input feature vector)

45 ^ Reply • Share >



disqus_u1549NhgT7 • 2 years ago

Guys, when trying to install (pip install rnn_theano) I get the following: No matching distribution found for rnn_theano... Any thoughts on how I should install this?



Anthony Caterini → disqus_u1549NhgT7 • 2 years ago

You shouldn't be trying to install rnn_theano from pip, since it's not a package (as far as I know). You should instead clone the repository from github.



Gabrer • 3 years ago

In "train_with_sdg" there is a typos in the "if" statement:

```
if (len(losses) > 1 and losses[-1][1] > losses[-2][1]):
```

insted of:



Patrick Stetz → Gabrer • a year ago

Am I not understanding a joke? Those are the same



Martin Kamp Dalgaard → Patrick Stetz • 9 months ago • edited

What it actually says is "if (len(losses) > 1 and losses[-1][1] > losses[-2][1]):". ">" is HTML entity for ">" (http://www.fileformat.info/...: -)

Edit: I realize now why the same signs appear in Gabrer's and my comments. It says "& g t;" (without spaces) in the blog (which I also tried to write), but it automatically becomes ">" when posting the comment.



Boussad • 2 years ago

Hi Denny,

Let me ask a question about your RNN implementation. It is really great but I think you implement a standard batch gradient descent backpropagation not SGD because you take all the samples from your dataset at every epoch. This is why it takes a very long time. In a pure SGD or online GD, only one sample is taken into account. This latter is much faster but may be inneficient given noisy data. This is why mini-batch SGD is generally used: for instance use 100 samples at every epoch --> 8000/100 epochs to scan the whole dataset (take 100 without replacement case). You save a lot of time then ...

Again, very great blog!

```
5 ^ | V • Reply • Share >
```



Mostafa Kotb • 3 years ago

i'm confused in calculating loss:

L += -1 * np.sum(np.log(correct_word_predictions))

where is y[i], it suppose to be:

L += -1 * np.sum(y[i] * np.log(correct_word_predictions))

am i missing something?

2 ^ V • Reply • Share >



easton YI A Mostafa Kotb • 2 years ago • edited

y[i] is the ith sentence

the ideal distribution of `correct_word_predictions` is the all-ones vector([1, 1, 1, 1, ...])

and the calculating loss here is equal to `L += -1 *

np.sum(1*np.log(correct_word_predictions))`

no questions here.

1 ^ V • Reply • Share >



Nitin Agarwal • 3 years ago

Thanks for such great post. I was facing problem in understanding the working of theano.scan(). It would be really very helpful if someone could explain how the parameters are being passed to forward prop step and output being generated.

2 ^ Reply • Share >



Chloe • 2 years ago

Hi, your post is really great!

I don't understand something, when the network create a sentence he takes words that have the biggest probability. How did he do to not create the same sentence every time?

1 ^ Reply • Share



Aman Sinha • 2 years ago

In defintion of generate sentence(model) funtion:

next word probs = model.forward propagation(new sentence)

might cause problem because it return two set of values.

It should be:

next word probs, s = model.forward propagation(new sentence)

1 ^ Reply • Share



Anthony Caterini • 2 years ago

Hi Denny,

Let me just say first of all that this is a great tutorial. I would just like to make a quick comment and ask a question about the new version of Theano - in particular the new apparray backend.

In your gradient_check_theano function, the variable backprop_gradient is now of type gpuarray. This means that estimated_gradient cannot be subtracted from it, throwing an error. To fix this, I added the following lines in:

bp_gpu = bptt_gradients[pidx][ix]
backprop_gradient = np.zeros(bp_gpu.shape, dtype=str(bp_gpu.dtype))
bp_gpu.read(backprop_gradient)

Now backprop_gradient is a numpy array of length 1, which allows it to be added to estimated_gradient. I wanted to ask something about this though: is this the new standard for extracting elements from a gpuarray, or is there an easier way to do this?

Thanks



Daniel Gaeta • 3 years ago • edited

I have a small question about the implementation of the loss function.

How is 'np.sum(np.log(correct_word_predictions))' mathematically equivalent to the equations y n * log(o n) in the loss function equation.

In particular, I am referring to the line:

L += -1 * np.sum(np.log(correct_word_predictions))

Update:

Am I correct in reasoning this is because y_n is 100% so we multiply by -1? That makes sense.



Daniel Seita → Daniel Gaeta • 2 years ago

It's because the targets are one-hot vectors.



Charles → Daniel Seita • 2 years ago

Am I crazy because I do not see where the data is put into one hot vectors?

• Reply • Share >

Show more replies



Kazi Nazmul Haque • 3 years ago

simple implementation of RNN and GRU in Tensorflow with scan and map ops. https://github.com/KnHuq/Dy...



Nanya → Kazi Nazmul Haque • 3 years ago

Hi,Kazi.Great repo.I'm new to RNN,and I am wondering about changing mnist data

in you ran implementation with my own cey files is that noscible?

Recurrent Neural Networks Tutorial, Part 2 - Implementing a RNN with Python, Numpy and Theano - WildML ווו איטע ווווו וווויףוכוווכווומוויטוו איווו וווא טאוו באי וווכא.וא נוומג איטאוווים:



jiang → Nanya • 2 years ago

你好,我现在也是在做CNN对excel进行分类,对比。我之前有对CSV文件分 类的例子。只是想问下你为什么要把excel文件转换成CSV而不是直接处理 呢?这里你遇到了那些问题?现在解决了吗?

∧ V • Reply • Share >

Show more replies



sodis • 3 years ago

sorry, I'm confused, in here:

def forward propagation(self, x):

The total number of time steps

T = len(x)

During forward propagation we save all hidden states in s because need them later.

We add one additional element for the initial hidden, which we set to 0

s = np.zeros((T + 1, self.hidden dim))

s[-1] = np.zeros(self.hidden dim)

why did s[-1] was assigned twice?

1 ^ V • Reply • Share >



Dhirendra → sodis • 2 years ago • edited

wrong post

∧ V • Reply • Share >



Evalds Urtans → sodis • 3 years ago

I think this is not needed:)

∧ V • Reply • Share >



Avatar This comment was deleted.



Avatar This comment was deleted.



Denny Britz Mod → Guest • 3 years ago

Good luck;)

∧ V • Reply • Share >

Show more replies



Mardan Hoshur • 3 years ago

Thank you for your great post.

I have one suggestion on coding. I understand that the parameter "x" in "def forward_propagation(self, x)" is a sentence (list of words). But the "x" in "def calculate_total_loss(self, x, y)" is a list of sentences. I was confused at the beginning because I got your idea largely by reading your code. It would be helpful for beginners like me if the second "x" (and "y") denoted with something else, such as capital "X".



Denny Britz Mod → Mardan Hoshur • 3 years ago

Thanks! I can see how that would be confusing. I'll try to fix it in the Github code.



Shad Akhtar • 3 years ago

Really nice post. It helped me a lot in understanding RNN. However, I have one doubt.

As far as I understand, the network has 8000(vacab_size) input units, means each word should be represented by a vector of 8000 values. But every instance of X_train i.e. X_train[i], has vocabulary index of all the words in a sentence. Now, during call of sgd_step(X_train[10],....) how exactly single index value is converted into vector of size 8000? or it is not?

I want to use RNN for sentiment analysis. In my case each word is represented by a vector length of 300. So, how should I provide input to sgd_step()? Should X_train[i] be populated as [300 x 10] for ten word sentence?

I'm still in the learning phase of deep learning. So, please correct me if I'm wrong.



Denny Britz Mod → Shad Akhtar • 3 years ago

The vocabulary size of 8000 means that we only consider the 8000 most common words, it's not related to the embedding dimension, which is what you are talking about. The embedding dimension above is the height of the U matrix, it's set to 100, but you can set it to 300 if you want to. U would then be a 300 x 8000 matrix, where each of the 8000 words is represented by a 300-dimensional vector.



Shad Akhtar → Denny Britz • 3 years ago

Thanks Denny. I missed a point earlier, now it is clear.



LI Bo • 3 years ago

Hi, this post perfectly solved all my questions about RNN, Great job!

However I found a problem when I was running the code. When I was trying the "gradient check" function, the result is different between RNNNumpy and RNNTheano. The "gradient check" function responses "Pass" when I was using the RNNNumpy class as you posted. But it failed at RNNTheano class.

Recurrent Neural Networks Tutorial, Part 2 – Implementing a RNN with Python, Numpy and Theano – WildML

np.random.seed(10)

grad check vocab size = 100

model = RNNTheano(grad_check_vocab_size, 10, hidden_dim = 100, bptt_truncate = 1000)

gradient_check_theano(model, [0,1,2,3], [1,2,3,4])

Above is the code when I was trying to "gradient check" using the RNNTheano class. As you can see, the parameters are the same as the "gradient check" function in RNNNumpy. But the result is:

Performing gradient check for parameter U with size 1000.

Gradient check ERROR: parameter = U ix = (0, 0)

+h Loss: 18.307188 -h Loss: 18.307188

Estimated_gradient: 0.000000
Backpropagation: -0.214507
Relative Error: 1.000000

Is that because the "T.grad" function is not efficient as BPTT algorithm or I did something wrong with the code?

Thanks!

1 ^ Reply • Share



Denny Britz Mod → LI Bo • 3 years ago

My guess is that there's something wrong with the code (or the grad check function). Since Theano calculates the gradients automatically it should definitely pass, and it did for me.

The fact that your estimated gradient is 0 and the calculated one isn't looks strange. Make sure your loss function is correct...



Botty Dimanov → Denny Britz • 3 years ago

Hey Denny, exquisite explanation on RNNs I am eternally grateful!

I have ran into the same issue, which is even more peculiar because contrary to any logic when the bptt_truncate parameter gets changed, the **estimated** gradient value follows suit.

e.g.:

grad_check_vocab_size = 100
np.random.seed(10)
model_check = RNNTheano(grad_check_vocab_size, 10, bptt_truncate=5)
gradient_check_theano(model_check,[0,1,2,3], [1,2,3,4])

output:

Performing gradient check for parameter U with size 1000.

Gradient check for parameter U passed.

Performing gradient check for parameter V with size 1000.

Gradient Check ERROR: parameter=V ix=(0, 0)

+h Loss: 18.307188 -h Loss: 18.307188

see more

Show more replies



Abhishek Shivkumar • 4 years ago

This is a wonderful blog post. Can't wait to see the next blog post on LSTM networks, something I always struggled to understand better. Also, can you recommend a good material/videolecture for understanding BPTT algorithm? Of course, I am sure you will be covering it in the next post - if you can supplement it with relevant papers/links, I am sure it would be great. Thanks again! Great job!



Denny Britz Mod → Abhishek Shivkumar • 4 years ago

Thanks Abhishek! This is a great post on LSTMs, it has really helped me understand them better: http://colah.github.io/post.... I probably won't be able to do better than that, but I'll supplement it with some code.

I can't think of any good video lectures on BPTT off the top of my head, but BPTT is really just a fancy name for applying standard backpropagation on the "unrolled" network and adding up the gradients. So I'd recommend really trying understand how standard backpropagation works and deriving the formula. Then BPTT will become clear as well.



Abhishek Shivkumar → Denny Britz • 4 years ago

Thanks for the reply Denny. I know the standard backprop algorithm and will go ahead and learn the BPTT. I have one question regarding your code. In the forward_prop_step method, why do you return the zeroeth element of o_t as return [o_t[0], s_t]? Shouldn't we return the element that has the maximum probability? I am trying to understand why o_t[0] from the output of the softmax function that outputs probability values?

Show more replies



Nguyễn Học • 9 months ago
Is this using Teacher Forcing?

∧ V • Reply • Share >

Blake Conrad • a year ago

Www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-2-implementing-a-language-model-rnn-with-python-numpy-and-theano/



Recurrent Neural Networks Tutorial, Part 2 – Implementing a RNN with Python, Numpy and Theano – WildML LACELLINI tutorial. I am study on your difficulties and snapes of the variable matrices. I would think that with an input vector of 1X8000, the following weight matrice should be 8000X8000, not 100X8000. That doesn't seem compatible. Could you explain this a bit further? Again, GREAT tutorial.



Omar S • a year ago

Hi thanks for this tutorial.

I am a bit lost as to the one-hot encoding. Where in the RNN did you implement it? because you said you won't do it in the pre-processing phase.

Thanks



Thuan • a year ago

How did you make the true labels "y", from which the cost function was calculated? Because in previous post, you said no label were needed, so how can we calculate the cost function.

Thanks



Parvez Khan • a year ago

It taking, too much time to train the model. Is there anything that we can do to reduce it???

• Reply • Share >



Arpit • a year ago

How can I download the .csv file to start with?



Jae Duk Seo • a year ago

Great post. Thank you!



舒展 • 2 years ago

Thanks for your great post. I got a question. How do you generate different setenses with same start token?



Akanksha Dwivedi • 2 years ago

Hi.

Can anyone please share the RNNTheano class code. I cannot seem to find the link. Thank you.