

Deep Learning: Homework 3

Name: Yifan Shen

SN: 15015762

Start date: 12th Feb 2018

Due date: 12nd March 2018, 11:55 pm

How to Submit

When you have completed the exercises and everything has finished running, click on 'File' in the menu-bar and then 'Download .ipynb'. This file must be submitted to Moodle named as **studentnumber_DL_hw3.ipynb** before the deadline above.

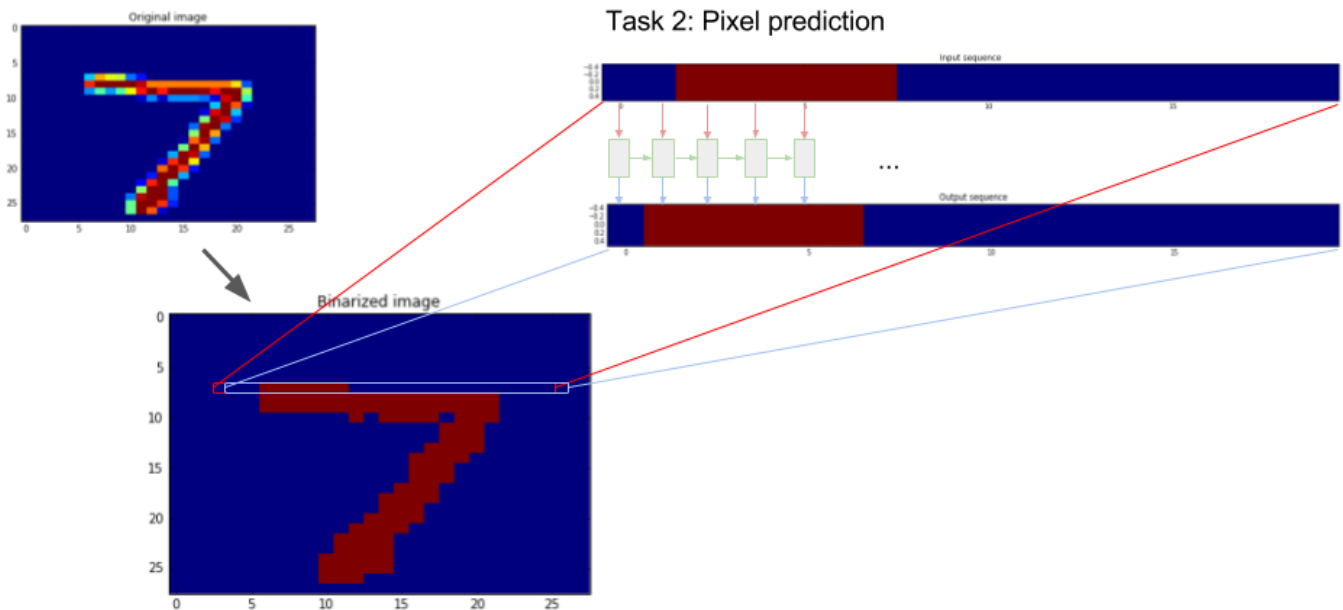
Also send a **sharable link** to the notebook at the following email: ucl.coursework.submit@gmail.com. You can also make it sharable via link to everyone, up to you.

Please compile all results (table in Q2) and all answers to the understanding/analysis results questions (Q1 and Q3), into a PDF. Name convention: **studentnumber_DL_hw3.pdf**. Do not include any of the code (we will use the notebook for that).

Page limit: 7 pg (w/o the bonus question).

MNIST as a sequence

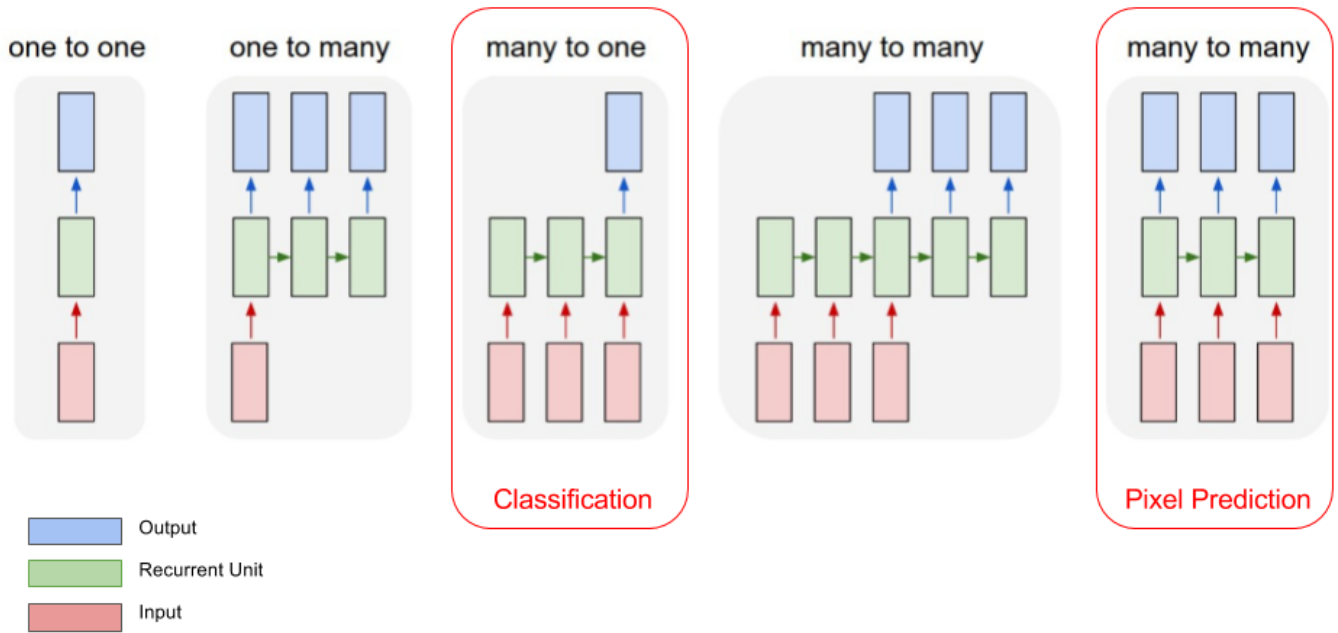
In this assignment we will be using the MNIST digit dataset (<https://yann.lecun.com/exdb/mnist/>). The dataset contains images of hand-written digits (0 – 9), and the corresponding labels. The images have a resolution of 28×28 pixels. This is the same dataset as in Assignment 1, but we will be using this data a bit differently this time around. Since this assignment will be focusing on recurrent networks that model sequential data, we will be looking at each image as a sequence: the networks you train will be "reading" the image one row at a time, from top to bottom (we could even do pixel-by-pixel, but in the interest of time we'll do row-by-row which is faster). Also, we will work with a binarized version of MNIST -- we constrain the values of the pixels to be either 0 or 1. You can do this by applying the method `binarize`, defined below, to the raw images.



- We take the MNIST images, binarise them, and interpret them as a sequence of pixels from top-left to bottom-right. ("Task 2" refers to the next homework, wherein you will be using the sequence for pixel prediction).

Recurrent Models for MNIST

As discussed in the lectures, there are various ways and tasks for which we can use recurrent models. A depiction of the most common scenarios is available in the Figure below. In this assignment and the following one we will look at two of these forms: **many-to-one** (sequence to label/decision) and the **many-to-many** scenario where the model receives an input and produces an output at every time step. You will use these to solve the following tasks: i) classification (*this homework*), ii) pixel prediction (*next homework*) and iii) in-painting (*next homework*).



- (Figure adapted from Karpathy's The Unreasonable Effectiveness of Recurrent Neural Networks (<http://karpathy.github.io/2015/05/21/rnn-effectiveness>)). You will be implementing variants of *many-to-one* for classification (in this homework), and *many-to-many* for prediction (in the next homework).

In [0]:

1

18## Q1: Understanding LSTM vs GRU (30 pts) Before going deeper into your practical tasks, take some time to revise and make sure you understand the two major types of recurrent cells you will be using in this assignment: Long-Short Term Memory Units (LSTM) first introduced by Hochreiter and Schmidhuber [1997] and the more recent Gated Recurrent Units (GRU) by Cho et al. [2014]. Once you have done this, answer the following questions:

Vanilla RNN:

$$h_t = f(W_{hx}x_t + W_{hh}h_{t-1})$$

LSTM:

$$\text{Input gate: } i_t = \sigma(W^{(i)}x_t + U^{(i)}h_{t-1})$$

$$\text{Forget gate: } f_t = \sigma(W^{(f)}x_t + U^{(f)}h_{t-1})$$

$$\text{Output gate: } o_t = \sigma(W^{(o)}x_t + U^{(o)}h_{t-1})$$

$$\text{Memory cell: } \tilde{c}_t = \tanh(W^{(c)}x_t + U^{(c)}h_{t-1})$$

$$\text{Final memory: } c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$\text{Hidden state: } h_t = o_t \circ \tanh(c_t)$$

GRU:

$$\text{Update gate: } z_t = \sigma(W^{(z)}x_t + U^{(z)}h_{t-1})$$

Reset gate: $r_t = \sigma(W^{(r)}x_t + U^{(r)}h_{t-1})$

Intermediate memory content: $\tilde{h}_t = \tanh(Wx_t + r_t \circ U h_{t-1})$

Hidden state: $h_t = z_t \circ h_{t-1} + (1 - z_t) \circ \tilde{h}_t$

1. Can LSTMs (and respectively GRUs) just store the current input in the state (c_t for LSTM and h_t for GRU, in the class notation) for the next step? If so, give the gates activation that would enable this behaviour. If not, explain why not. [10 pts]

Answer 1: Both LSTM and GRU can do this. For LSTM, $i_t > 0$ and $f_t = 0$. For GRU, $z_t = 1$ and $r_t = 0$

2. Can LSTMs (and respectively GRUs) just store a previous state into the current state and ignore the current input? If so, give the gates' activation that would enable this. If not, explain why not. [10 pts]

Answer 2: Both LSTM and GRU can do this. For LSTM, $i_t = 0$ and $f_t > 0$. For GRU, $z_t = 0$ and $r_t > 0$

3. Are GRUs a special case of LSTMs? If so, give the expression of the GRU gates in term of LSTM's gates (o_t, i_t, f_t). If not, give a counter-example. Assume here the same input. [10 pts]

Answer 3: GRUs are not a special case of LSTMs.

We propose a special case of GRU that cannot be attained by LSTM.

Let $z_i = 1, \forall i \in 1, 2, \dots, t$ then we have $h_t = h_{t-1} = h_{t-2} = \dots = h_1$. If we want match the same result from LSTM side, we need $h_t = o_t \circ \tanh(f_t \circ c_{t-1} + i_t \circ \tanh(W^{(c)}x_t + U^{(c)}h_{t-1})) = h_{t-1}$. Thus, in theory it could be possible for the LSTM to attain this result $h_t = h_{t-1} = \dots = h_1$ by tuning the weights of the two tanh layers so that they operate an identity transformation, and by forgetting everything that is in the cell state and replacing it with the last output. Though, such tuning would be valid only for a specific value of the output. Therefore, in general, the LSTM cannot decide, at some point, to start reproducing the last output.

Word limit: 1000 words or less

Q2: Implementation. Line-by-Line MNIST Classification (50 pts)

In this part you will train a number of many-to-one recurrent models that takes as input: an image (or part of an image) as a sequence (row by row) and after the last input row produces, as output, a probability distribution over the 10 possible labels (0 – 9). The models will be trained using a cross-entropy loss function over these output probabilities.

Optimization

Use the Adam optimizer (with default settings other than the learning rate) for training.

[Optional] Sometimes dropout has been shown to be beneficial in training recurrent models, so feel free to use it or any other form of regularization that seems to improve performance. It might be also worth trying out batch-normalization. [Reference \(https://arxiv.org/pdf/1603.09025.pdf\)](https://arxiv.org/pdf/1603.09025.pdf).

Models: Your models will have the following structure:

1. [(Red Block)] The *input* (current binarised row of pixels) can be fed directly into the recurrent connection without much further pre-processing. The only thing you need to do is have an affine transformation to match the dimensionality of the recurrent unit, i.e. one of (32, 64, 128).

2. [(Blue Block)] The *output* (probabilities over the 10 classes) is produced by looking at the last output of the recurrent units, transforming them via an affine transformation.
3. [(Green Block)] For the *recurrent* part of the network, please implement and compare the following architectures:
 - LSTM with 32, 64, 128 units. [15 pts]
 - GRU with 32, 64, 128 units. [15 pts]
 - stacked LSTM: 3 recurrent layers with 32 units each. [10 pts]
 - stacked GRU: 3 recurrent layers with 32 units each. [10 pts]

Your network should look like:

Input \Rightarrow RNN cell \Rightarrow Relu \Rightarrow Fully connected \Rightarrow Relu \Rightarrow Fully connected \Rightarrow Output

You might find the function `tf.nn.dynamic_rnn` useful.

Hyper-parameters

For all cases train the model with these hyper-parameter settings:

- `num_epochs=10, learning_rate=0.001, batch_size=256, fully_connected_hidden_units=64`

With these hyper-parameters you should be comfortably above 95% test set accuracy on all tasks. (Feel free to try other settings, there are certainly better choices, but please report the results with these exact hyper-parameters). Please report the *cross-entropy* and the *classification accuracy* for the *test set* of the models trained. Use the `plot_summary_table` method below to format the table.

##Q3: Analyse the results (20 pts + 10 pts)

1. How does this compare with the results you obtained in the first assignment(DL1), when training a model that "sees" the entire image at once? Explain differences. [5 pts]

Answer: The results found by LSTM and GRU are slightly better than the normal neural networks (first three settings in the first assignment) but slightly worse than CNN.

Reason: Normal neural networks "sees" the entire image at once while RNNs "sees" a row of the image each time. Therefore RNNs captures more information by "looking" the image step by step (a row per time, for 28 times), resulting better performance than normal neural networks. For CNNs, it "sees" a 3x3 window of the image everytime, and all pixels have been "seen" for several times. Therefore CNNs stores much even more information than RNNs, resulting the best performance.

1. Let us take a look closer look at one of the trained models: say GRU (32). Plot the outputs of the RNN layer and hidden state over time. In particular, look at the first 3-5 time steps. Plot the input image along side. You can use python `plt.imshow(output_GRU_over_time)` for these, where `output_GRU_over_time.shape` is (T=28,hidden_units) dimensional. What do you observe? Show at least one pair of these plots to support your observation(s). [5 pts]

Answer: In general, the stripes of "output_GRU_over_time" plot is vague at the very beginning, but becomes clear with time steps. At the very beginning, the input are all 0s (have not reach the true number yet), which are meaningless. Therefore, it is nor surprising that the model cannot distinguish the digit at the beginning and resulting very small value in "output_GRU_over_time".

1. Now, look at the last 3-5 time steps. What do you observe? When is the classification decision made? To validate your answer to the second part of the question, provide the classification predictions for the last 5 time steps -- that is, pretend whatever the output of the GRU is at that time step is in fact the last output in the computation, and feed that into the classification mapping. [10 pts]

Answer: The stripes of "output_GRU_over_time" of the last 3-5 times becomes very clear, which means the classification has been made. In the final part of my code, I have printed out the decision made each time. In this case, our label is 5, we can see that the decision is made at the 18th step. I have tried a couple of different examples, the decisions are usually made at the later steps (around 16 to 25)

Word limit: 500 words or less

[Bonus] Let's looking inside the computation. Take one of the previous LSTM models (for simplicity pick one of the one-layer recurrence models) and track the status of the gates (o_t, i_t, f_t) over time. Note that if you used the provided RNNCell wrappers (BasicLSTMCell and co) in tensorflow, these will keep this information hidden, so you will need to implement your own version of this recurrence layer, or mirror the one in tensorflow, but now exposing these hidden variables. A bit of warning: this is not trivial and will require some thinking on the coding side, but it will also provide you with a more informative way of visualizing the inner computation. [10 pts]

Word limit: 300 words or less

Imports and utility functions (do not modify!)

In [0]:

```

1 # Import useful libraries.
2 import tensorflow as tf
3 import matplotlib.pyplot as plt
4 from tensorflow.examples.tutorials.mnist import input_data
5 import numpy as np
6
7 # Binarize the images
8 def binarize(images, threshold=0.1):
9     return (threshold < images).astype('float32')
10
11 # Import dataset with one-hot encoding of the class labels.
12 def get_data():
13     return input_data.read_data_sets("MNIST_data/", one_hot=True)
14
15 # Placeholders to feed train and test data into the graph.
16 # Since batch dimension is 'None', we can reuse them both for train and eval.
17 def get_placeholders():
18     x = tf.placeholder(tf.float32, [None, 784])
19     y = tf.placeholder(tf.float32, [None, 10])
20     return x, y
21
22 # Generate summary table of results. This function expects a dict with the
23 # following structure: keys of 'LSTM' or 'GRU' and the values for each key are
24 # list of tuples consisting of (test_loss, test_accuracy), and the list is
25 # ordered as the results from 32 units, 64 units, 128 units, 3 x 32 units, i.e.
26 # {
27 #   'LSTM': [(loss,acc), (loss, acc), (loss, acc), (loss, acc)]
28 #   'GRU': [(loss,acc), (loss, acc), (loss, acc), (loss, acc)]
29 # }
30 def plot_summary_table(experiment_results):
31     # Fill Data.
32     cell_text = []
33     columns = ['(1 layer, 32 units)', '(1 layer, 64 units)', '(1 layer, 128 units)',
34               '(3 layers, 32 units)']
35     for k, v in experiment_results.items():
36         rows = ['Test loss', 'Test accuracy']
37         cell_text=[[],[]]
38         for (l, _) in v:
39             cell_text[0].append(str(l))
40         for (_, a) in v:
41             cell_text[1].append(str(a))
42
43     fig=plt.figure(frameon=False)
44     ax = plt.gca()
45     the_table = ax.table(
46         cellText=cell_text,
47         rowLabels=rows,
48         colLabels=columns,
49         loc='center')
50     the_table.scale(2, 8)
51     # Prettify.
52     ax.patch.set_facecolor('None')
53     ax.xaxis.set_visible(False)
54     ax.yaxis.set_visible(False)
55     ax.text(-0.73, 0.9, k, fontsize=18)

```

Train Models

Generate summary table of results. This function expects a dict with the following structure: keys of 'LSTM' or 'GRU' and the values for each key are a list of tuples consisting of (test_loss, test_accuracy), and the list is ordered as the results from 32 units, 64 units, 128 units, 3 x 32 units, i.e. expected dictionary (final performance only):

```
{
  'LSTM': [(loss, acc), (loss, acc), (loss, acc), (loss, acc)]
  'GRU': [(loss, acc), (loss, acc), (loss, acc), (loss, acc)]
}
```


In [0]:

```

1  # Your models here
2
3  num_epochs = 10
4  learning_rate = 0.001
5  batch_size = 256
6  fully_connected_hidden_units = 64
7  num_classes = 10
8
9  initializer = tf.contrib.layers.xavier_initializer()
10
11 def RNN(RNN_type, num_hidden_unit, num_layers):
12     time_steps = 28
13     num_units = 28
14     tf.reset_default_graph()
15
16     mnist = get_data()
17     eval_mnist = get_data()
18     x, y_ = get_placeholders()
19     x_reshape = tf.reshape(x, [-1, 28, 28])
20
21     if RNN_type == 'LSTM':
22         RNN_layers = [tf.contrib.rnn.BasicLSTMCell(size) for size in \
23                        ([num_hidden_unit] * num_layers)]
24     elif RNN_type == 'GRU':
25         RNN_layers = [tf.contrib.rnn.GRUCell(size) for size in \
26                        ([num_hidden_unit] * num_layers)]
27
28     cell = tf.contrib.rnn.MultiRNNCell(RNN_layers)
29
30     rnn_outputs, state = tf.nn.dynamic_rnn(cell, x_reshape, dtype=tf.float32) \
31     # rnn_outputs: batch_size * 28 * num_hidden
32     # choose the last one by using -1 (many to one)
33     rnn_outputs_unstack = tf.unstack(rnn_outputs, time_steps, 1)[-1] \
34     # rnn_outputs_unstack: batch_size * num_hidden
35
36     W_1 = tf.Variable(initializer([num_hidden, fully_connected_hidden_units]))
37     b_1 = tf.Variable(initializer([fully_connected_hidden_units]))
38     hid_1 = tf.nn.relu(tf.matmul(rnn_outputs_unstack, W_1) + b_1)
39
40     W_2 = tf.Variable(initializer([fully_connected_hidden_units, \
41                                   fully_connected_hidden_units]))
42     b_2 = tf.Variable(initializer([fully_connected_hidden_units]))
43     hid_2 = tf.nn.relu(tf.matmul(hid_1, W_2) + b_2)
44
45     W_3 = tf.Variable(initializer([fully_connected_hidden_units, num_classes]))
46     b_3 = tf.Variable(initializer([num_classes]))
47     y = tf.matmul(hid_2, W_3) + b_3
48
49     # loss
50     cross_entropy = \
51     tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(labels = y_, logits = y))
52
53     # update
54     train_step = tf.train.AdamOptimizer(learning_rate).minimize(cross_entropy)
55
56     # evaluation metric
57     correct_prediction = tf.equal(tf.argmax(tf.nn.softmax(y), 1), tf.argmax(y_, 1))
58     accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
59

```

```
60 i, train_accuracy, test_accuracy = 0, [], []
61
62 with tf.train.MonitoredSession() as sess:
63
64     while mnist.train.epochs_completed < num_epochs:
65
66         i += 1
67         batch_xs, batch_ys = mnist.train.next_batch(batch_size)
68
69         sess.run(train_step, feed_dict = {x: binarize(batch_xs), y_: batch_ys})
70
71         if i % (100) == 0:
72
73             test_acc = accuracy.eval(session = sess, feed_dict= \
74                 {x: binarize(mnist.test.images), y_: mnist.test.labels})
75             loss = sess.run(cross_entropy, \
76                 feed_dict = {x: binarize(batch_xs), y_: batch_ys})
77             print('%d th test acc: %s loss: %s' %(i, test_acc, loss))
78
79     return (loss, test_acc)
```

In [41]:

```

1 settings = [\
2     ['LSTM',32,1],\
3     ['LSTM',64,1],\
4     ['LSTM',128,1],\
5     ['LSTM',32,3],\
6     ['GRU',32,1],\
7     ['GRU',64,1],\
8     ['GRU',128,1],\
9     ['GRU',32,3]]
10
11 experiment_results = {'LSTM':[], 'GRU':[]}
12
13 for RNN_type, num_hidden, num_layers in settings:
14
15     print(RNN_type, num_hidden, num_layers)
16
17     experiment_results[RNN_type].append(RNN(RNN_type, num_hidden, num_layers))
18
19 print(experiment_results)
20
21

```

LSTM 32 1

Extracting MNIST_data/train-images-idx3-ubyte.gz

Extracting MNIST_data/train-labels-idx1-ubyte.gz

Extracting MNIST_data/t10k-images-idx3-ubyte.gz

Extracting MNIST_data/t10k-labels-idx1-ubyte.gz

Extracting MNIST_data/train-images-idx3-ubyte.gz

Extracting MNIST_data/train-labels-idx1-ubyte.gz

Extracting MNIST_data/t10k-images-idx3-ubyte.gz

Extracting MNIST_data/t10k-labels-idx1-ubyte.gz

INFO:tensorflow:Graph was finalized.

INFO:tensorflow:Running local_init_op.

INFO:tensorflow:Done running local_init_op.

100 th test acc: 0.5623 loss: 1.2627294

200 th test acc: 0.7923 loss: 0.60181373

300 th test acc: 0.854 loss: 0.45864093

400 th test acc: 0.8835 loss: 0.44273466

500 th test acc: 0.9014 loss: 0.28687698

600 th test acc: 0.9156 loss: 0.29344687

700 th test acc: 0.9226 loss: 0.2072212

800 th test acc: 0.9256 loss: 0.15604701

#Results

In [42]:

```
1 plot_summary_table(experiment_results)
```

LSTM	(1 layer, 32 units)	(1 layer, 64 units)	(1 layer, 128 units)	(3 layers, 32 units)
Test loss	0.12693101	0.09858341	0.074853815	0.0649623
Test accuracy	0.9634	0.9772	0.9765	0.9709

GRU	(1 layer, 32 units)	(1 layer, 64 units)	(1 layer, 128 units)	(3 layers, 32 units)
Test loss	0.09399286	0.027650137	0.017552938	0.05128681
Test accuracy	0.9696	0.9781	0.9837	0.9766

Analysis of results

In [41]:

```

1
2 num_epochs = 10
3 learning_rate = 0.001
4 batch_size = 256
5 fully_connected_hidden_units = 64
6 num_classes = 10
7
8 initializer = tf.contrib.layers.xavier_initializer()
9
10 num_hidden_unit = 64
11 num_layers = 1
12
13 RNN_type = 'GRU'
14
15 time_steps = 28
16 num_units = 28
17 tf.reset_default_graph()
18
19 mnist = get_data()
20 eval_mnist = get_data()
21 x, y_ = get_placeholders()
22 x_reshape = tf.reshape(x, [-1, 28, 28])
23
24 if RNN_type == 'LSTM':
25     RNN_layers = [tf.contrib.rnn.BasicLSTMCell(size) for size in \
26                   ([num_hidden_unit] * num_layers)]
27 elif RNN_type == 'GRU':
28     RNN_layers = [tf.contrib.rnn.GRUCell(size) for size in \
29                   ([num_hidden_unit] * num_layers)]
30
31 cell = tf.contrib.rnn.MultiRNNCell(RNN_layers)
32
33 rnn_outputs, state = tf.nn.dynamic_rnn(cell, x_reshape, dtype=tf.float32) \
34 # rnn_outputs: batch_size * 28 * num_hidden
35 # choose the last one by using -1 (many to one)
36 rnn_outputs_unstack = tf.unstack(rnn_outputs, time_steps, 1)[-1]\
37 # rnn_outputs_unstack: batch_size * num_hidden
38
39 W_1 = tf.Variable(initializer([num_hidden_unit, fully_connected_hidden_units]))
40 b_1 = tf.Variable(initializer([fully_connected_hidden_units]))
41 hid_1 = tf.nn.relu(tf.matmul(rnn_outputs_unstack, W_1) + b_1)
42
43 W_2 = tf.Variable(initializer([fully_connected_hidden_units, \
44                               fully_connected_hidden_units]))
45 b_2 = tf.Variable(initializer([fully_connected_hidden_units]))
46 hid_2 = tf.nn.relu(tf.matmul(hid_1, W_2) + b_2)
47
48 W_3 = tf.Variable(initializer([fully_connected_hidden_units, num_classes]))
49 b_3 = tf.Variable(initializer([num_classes]))
50 y = tf.matmul(hid_2, W_3) + b_3
51
52 # loss
53 cross_entropy = \
54 tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(labels = y_, logits = y))
55
56 # update
57 train_step = tf.train.AdamOptimizer(learning_rate).minimize(cross_entropy)
58
59 # evaluation metric

```

```

60 correct_prediction = tf.equal(tf.argmax(tf.nn.softmax(y),1), tf.argmax(y_,1))
61 accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
62
63 i, train_accuracy, test_accuracy = 0, [], []
64
65 with tf.train.MonitoredSession() as sess:
66
67     while mnist.train.epochs_completed < num_epochs:
68
69         i += 1
70
71         batch_xs, batch_ys = mnist.train.next_batch(batch_size)
72
73         sess.run(train_step, feed_dict = {x: binarize(batch_xs), y_: batch_ys})
74
75         if i % (2149) == 0: # the last training step
76
77             test_acc = accuracy.eval(session = sess, \
78                                     feed_dict= {x: binarize(mnist.test.images), y_: mnist.test.labels})
79             loss = sess.run(cross_entropy, \
80                             feed_dict = {x: binarize(batch_xs), y_: batch_ys})
81             print('%d th test acc: %s loss: %s' % (i, test_acc, loss))
82             rnn_outputs_run = sess.run(rnn_outputs, feed_dict = \
83                                         {x: binarize(batch_xs), y_: batch_ys})
84             output_GRU_over_time = rnn_outputs_run[1,:,:]
85             binarized_input_image = binarize(batch_xs[1]).reshape(28,28)
86
87             # feed "output_GRU_over_time" to the trained model:
88
89             for j in range(output_GRU_over_time.shape[0]): # loop through time steps
90
91                 decision_result = sess.run(y, \
92                                             feed_dict = {rnn_outputs_unstack: output_GRU_over_time[j].reshape(1,-1)})
93                 decision = np.argmax(decision_result)
94
95                 print(j + 1, 'th time step, decision result is', decision)
96

```

```

Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz

```

INFO:tensorflow:Graph was finalized.

INFO:tensorflow:Running local_init_op.

INFO:tensorflow:Done running local_init_op.

2149 th test acc: 0.9779 loss: 0.10354345

```

1 th time step, decision result is 5
2 th time step, decision result is 5
3 th time step, decision result is 5
4 th time step, decision result is 5
5 th time step, decision result is 2
6 th time step, decision result is 9
7 th time step, decision result is 7
8 th time step, decision result is 7
9 th time step, decision result is 7
10 th time step, decision result is 7
11 th time step, decision result is 7

```

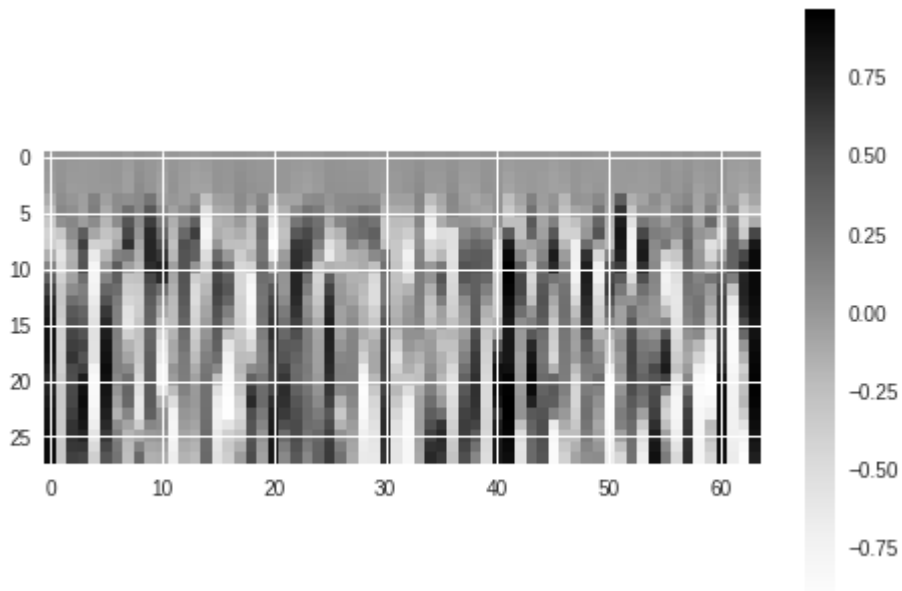
```
12 th time step, decision result is 0
13 th time step, decision result is 0
14 th time step, decision result is 0
15 th time step, decision result is 0
16 th time step, decision result is 0
17 th time step, decision result is 0
18 th time step, decision result is 5
19 th time step, decision result is 5
20 th time step, decision result is 5
21 th time step, decision result is 5
22 th time step, decision result is 5
23 th time step, decision result is 5
24 th time step, decision result is 5
25 th time step, decision result is 5
26 th time step, decision result is 5
27 th time step, decision result is 5
28 th time step, decision result is 5
```

In [42]:

```
1 plt.imshow(output_GRU_over_time)
2 plt.colorbar()
```

Out[42]:

<matplotlib.colorbar.Colorbar at 0x7f27893bd588>

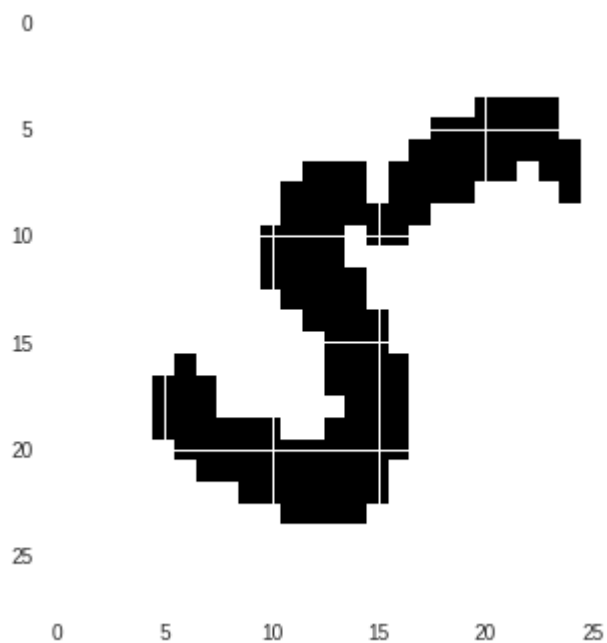


In [43]:

```
1 plt.imshow(binarized_input_image)
```

Out[43]:

<matplotlib.image.AxesImage at 0x7f2789dec400>

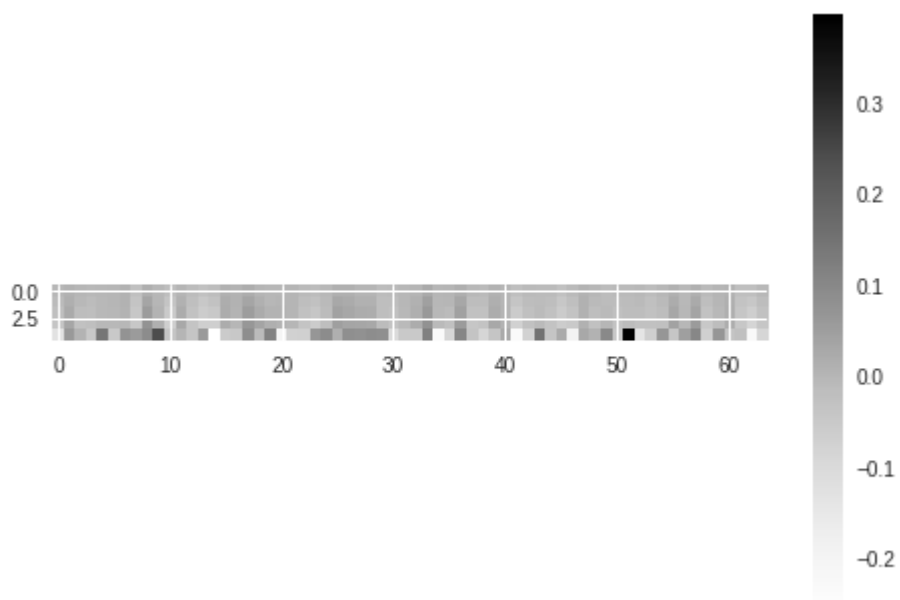


In [44]:

```
1 plt.imshow(output_GRU_over_time[:5,:])  
2 plt.colorbar()
```

Out[44]:

<matplotlib.colorbar.Colorbar at 0x7f2789f52dd8>

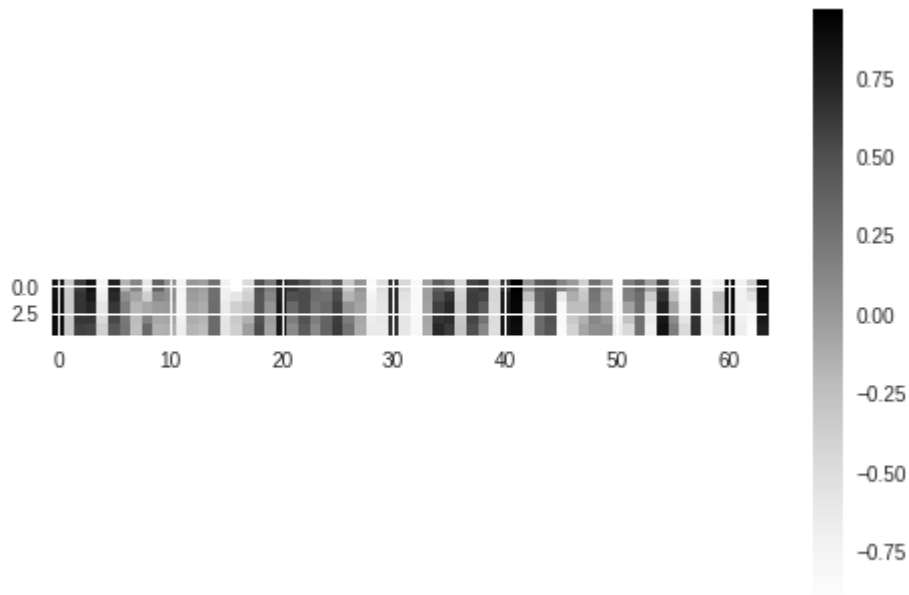


In [45]:

```
1 plt.imshow(output_GRU_over_time[-5:,:])  
2 plt.colorbar()
```

Out[45]:

<matplotlib.colorbar.Colorbar at 0x7f27899c8470>



Bonus Question:

In [0]:

```

1
2 class LSTM_cell(object):
3
4     """
5     LSTM cell object which takes 3 arguments for initialization.
6     input_size = Input Vector size
7     hidden_layer_size = Hidden layer size
8     target_size = Output vector size
9
10    """
11
12    def __init__(self, input_size, hidden_layer_size , x_batch_input):
13
14        # Initialization of given values
15        self.input_size = input_size
16        self.hidden_layer_size = hidden_layer_size
17        self.x_batch_input = x_batch_input
18
19        self.processed_input = self.x_batch_input
20
21        # Weights and Bias for input and hidden tensor
22        self.Wi = tf.Variable(tf.zeros(
23            [self.input_size, self.hidden_layer_size]))
24        self.Ui = tf.Variable(tf.zeros(
25            [self.hidden_layer_size, self.hidden_layer_size]))
26        self.bi = tf.Variable(tf.zeros([self.hidden_layer_size]))
27
28
29        self.Wf = tf.Variable(tf.zeros(
30            [self.input_size, self.hidden_layer_size]))
31        self.Uf = tf.Variable(tf.zeros(
32            [self.hidden_layer_size, self.hidden_layer_size]))
33        self.bf = tf.Variable(tf.zeros([self.hidden_layer_size]))
34
35
36        self.Wog = tf.Variable(tf.zeros(
37            [self.input_size, self.hidden_layer_size]))
38        self.Uog = tf.Variable(tf.zeros(
39            [self.hidden_layer_size, self.hidden_layer_size]))
40        self.bog = tf.Variable(tf.zeros([self.hidden_layer_size]))
41
42
43        self.Wc = tf.Variable(tf.zeros(
44            [self.input_size, self.hidden_layer_size]))
45        self.Uc = tf.Variable(tf.zeros(
46            [self.hidden_layer_size, self.hidden_layer_size]))
47        self.bc = tf.Variable(tf.zeros([self.hidden_layer_size]))
48
49        '''
50        Initial hidden state's shape is [1,self.hidden_layer_size]
51        In First time stamp, we are doing dot product with weights to
52        get the shape of [batch_size, self.hidden_layer_size].
53        For this dot product tensorflow use broadcasting. But during
54        Back propagation a low level error occurs.
55        So to solve the problem it was needed to initialize initial
56        hidden state of size [batch_size, self.hidden_layer_size].
57        So here is a little hack !!!! Getting the same shaped
58        initial hidden state of zeros.
59        '''

```

```

60
61     self.initial_hidden = self.x_batch_input[0, :, :]
62     self.initial_hidden= tf.matmul(
63         self.initial_hidden, tf.zeros([input_size, hidden_layer_size]))
64
65
66     self.initial_hidden=tf.stack([self.initial_hidden,self.initial_hidden,\
67                                   self.initial_hidden,self.initial_hidden,\
68                                   self.initial_hidden])
69
70 # Function for LSTM cell.
71 def Lstm(self, previous_hidden_memory_tuple, x):
72     """
73     This function takes previous hidden state and memory tuple with input a
74     outputs current hidden state.
75     """
76
77     previous_hidden_state, c_prev,_,_,_=tf.unstack(previous_hidden_memory_t
78
79     #Input Gate
80     i= tf.sigmoid(
81         tf.matmul(x,self.Wi)+tf.matmul(previous_hidden_state,self.Ui) \
82         + self.bi
83     )
84
85     #Forget Gate
86     f= tf.sigmoid(
87         tf.matmul(x,self.Wf)+tf.matmul(previous_hidden_state,self.Uf) \
88         + self.bf
89     )
90
91     #Output Gate
92     o= tf.sigmoid(
93         tf.matmul(x,self.Wog)+tf.matmul(previous_hidden_state,self.Uog) \
94         + self.bog
95     )
96
97     #New Memory Cell
98     c_ = tf.nn.tanh(
99         tf.matmul(x,self.Wc)+tf.matmul(previous_hidden_state,self.Uc) \
100         + self.bc
101     )
102
103     #Final Memory cell
104     c= f*c_prev + i*c_
105
106     #Current Hidden state
107     current_hidden_state = o*tf.nn.tanh(c)
108
109     return tf.stack([current_hidden_state,c,i,f,o]) # stack everything toge
110
111 # Function for getting all hidden state.
112 def get_states(self):
113     """
114     Iterates through time/ sequence to get all hidden state
115     """
116     # Getting all hidden state throuh time
117     all_hidden_states = tf.scan(self.Lstm,
118                                   self.x_batch_input,
119                                   initializer=self.initial_hidden,
120                                   name='states')

```

```
121
122         return all_hidden_states
123
124 # Function to convert batch input data to use scan ops of tensorflow.
125 def process_batch_input_for_RNN(batch_input):
126     """
127     Process tensor of size [5,3,2] to [3,5,2]
128     """
129     batch_input_ = tf.transpose(batch_input, perm=[2, 0, 1])
130     X = tf.transpose(batch_input_)
131
132     return X
133
```

In [60]:

```

1 num_epochs = 10
2 learning_rate = 0.001
3 batch_size = 256
4 fully_connected_hidden_units = 64
5 num_classes = 10
6
7 initializer = tf.contrib.layers.xavier_initializer()
8
9 num_hidden_unit = 32
10 num_layers = 1
11
12 time_steps = 28
13 num_units = 28
14 tf.reset_default_graph()
15
16 mnist = get_data()
17 eval_mnist = get_data()
18 x, y_ = get_placeholders()
19 x_reshape = tf.reshape(x, [-1, 28, 28])
20 x_processed = process_batch_input_for_RNN(x_reshape)
21
22 single_lstm_cell = LSTM_cell(num_units, 32, x_processed)
23
24 rnn_outputs_full = single_lstm_cell.get_states() # rnn_outputs_full: 28 * 5 * 2.
25
26 rnn_outputs = rnn_outputs_full[:,0,:,:] # rnn_outputs: 28 * 256 * 32
27
28 rnn_outputs_unstack = rnn_outputs[-1] # rnn_outputs_unstack: 256 * 32
29
30 W_1 = tf.Variable(initializer([num_hidden_unit, fully_connected_hidden_units]))
31 b_1 = tf.Variable(initializer([fully_connected_hidden_units]))
32 hid_1 = tf.nn.relu(tf.matmul(rnn_outputs_unstack, W_1) + b_1)
33
34 W_2 = tf.Variable(initializer([fully_connected_hidden_units, \
35                               fully_connected_hidden_units]))
36 b_2 = tf.Variable(initializer([fully_connected_hidden_units]))
37 hid_2 = tf.nn.relu(tf.matmul(hid_1, W_2) + b_2)
38
39 W_3 = tf.Variable(initializer([fully_connected_hidden_units, num_classes]))
40 b_3 = tf.Variable(initializer([num_classes]))
41 y = tf.matmul(hid_2, W_3) + b_3
42
43 # loss
44 cross_entropy = \
45 tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(labels = y_, logits = y))
46
47 # update
48 train_step = tf.train.AdamOptimizer(learning_rate).minimize(cross_entropy)
49
50 # evaluation metric
51 correct_prediction = tf.equal(tf.argmax(tf.nn.softmax(y),1), tf.argmax(y_,1))
52 accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
53
54 i, train_accuracy, test_accuracy = 0, [], []
55
56 with tf.train.MonitoredSession() as sess:
57
58     while mnist.train.epochs_completed < num_epochs:
59

```

```

60     i += 1
61
62     batch_xs, batch_ys = mnist.train.next_batch(batch_size)
63
64     sess.run(train_step, feed_dict = {x: binarize(batch_xs), y_: batch_ys})
65
66     if i % (2149) == 0: # the last training step
67
68         test_acc = accuracy.eval(session = sess, \
69             feed_dict= {x: binarize(mnist.test.images), y_: mnist.test.labels})
70         loss = sess.run(cross_entropy, \
71             feed_dict = {x: binarize(batch_xs), y_: batch_ys})
72         print('%d th test acc: %s loss: %s' %(i, test_acc, loss))
73         rnn_outputs_full_run = sess.run(rnn_outputs_full, \
74             feed_dict = {x: binarize(batch_xs), y_: batch_ys})
75         which_image = 0
76         binarized_input_image = binarize(batch_xs[which_image]).reshape(28,28)
77         lstm_i = rnn_outputs_full_run[:,2,which_image,:]
78         lstm_f = rnn_outputs_full_run[:,3,which_image,:]
79         lstm_o = rnn_outputs_full_run[:,4,which_image,:]
80

```

```

Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
Extracting MNIST_data/train-images-idx3-ubyte.gz
Extracting MNIST_data/train-labels-idx1-ubyte.gz
Extracting MNIST_data/t10k-images-idx3-ubyte.gz
Extracting MNIST_data/t10k-labels-idx1-ubyte.gz
INFO:tensorflow:Graph was finalized.
INFO:tensorflow:Running local_init_op.
INFO:tensorflow:Done running local_init_op.
2149 th test acc: 0.9515 loss: 0.12104231

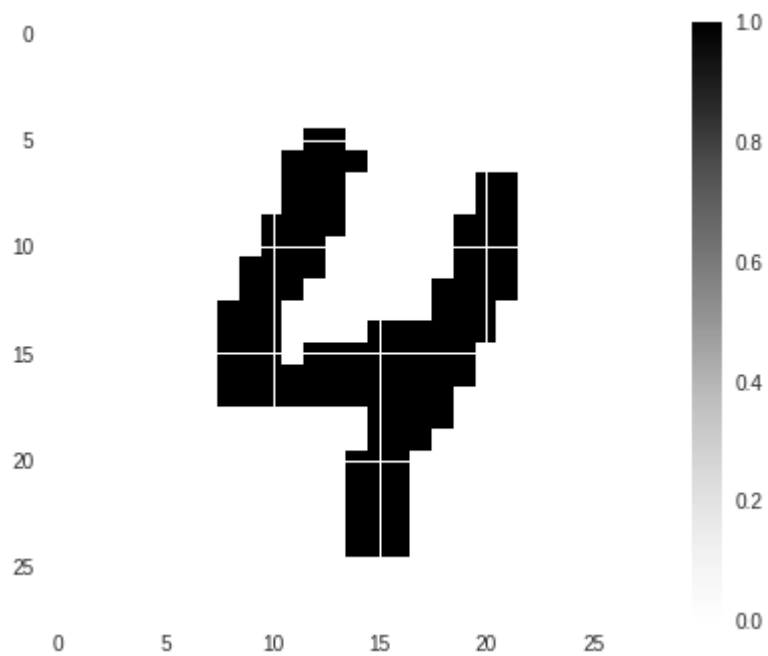
```

In [61]:

```
1 plt.imshow(binarized_input_image)
2 plt.colorbar()
```

Out[61]:

<matplotlib.colorbar.Colorbar at 0x7fe5f2936518>

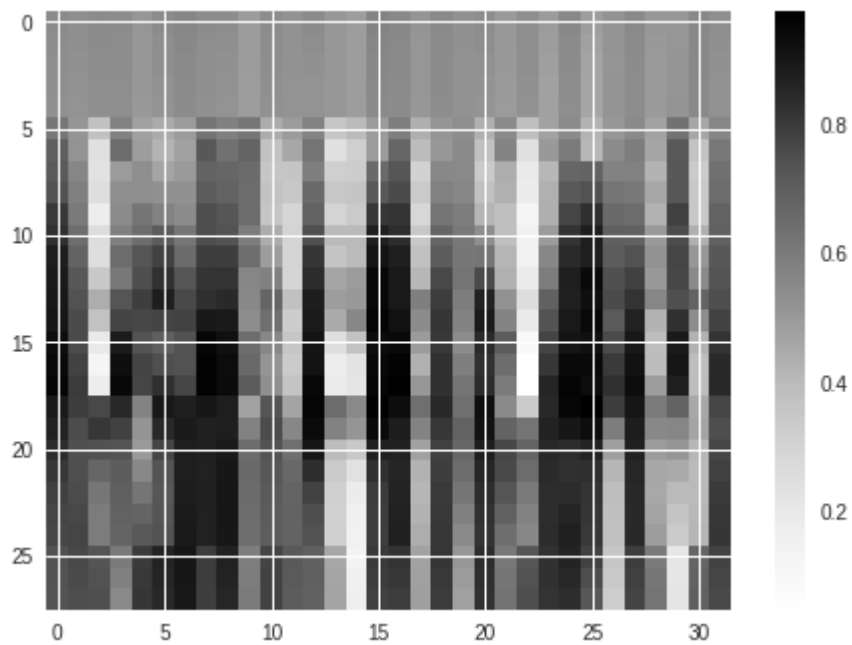


In [62]:

```
1 plt.imshow(lstm_i)
2 plt.colorbar()
```

Out[62]:

<matplotlib.colorbar.Colorbar at 0x7fe5f2f7cda0>

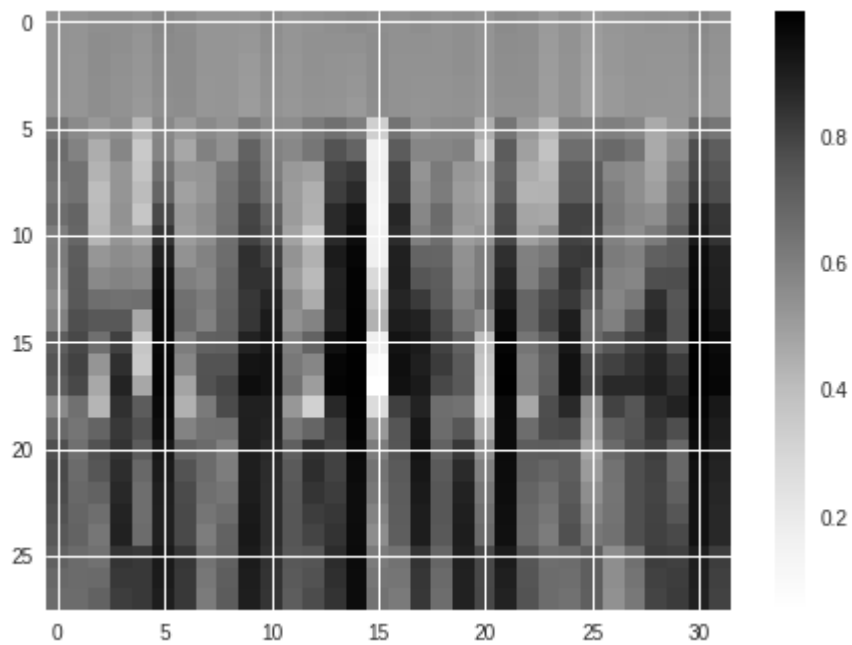


In [63]:

```
1 plt.imshow(lstm_f)
2 plt.colorbar()
```

Out[63]:

<matplotlib.colorbar.Colorbar at 0x7fe5f24dc9b0>

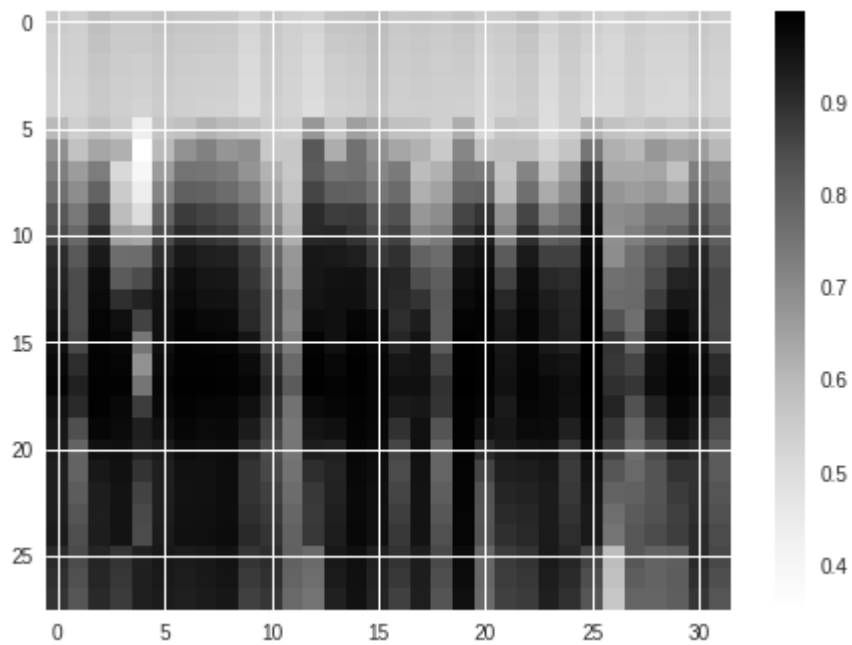


In [64]:

```
1 plt.imshow(lstm_o)
2 plt.colorbar()
```

Out[64]:

<matplotlib.colorbar.Colorbar at 0x7fe5f5ad4668>



In [0]:

```
1
```