

Image Captioning with LSTMs

In the previous exercise you implemented a vanilla RNN and applied it to image captioning. In this notebook you will implement the LSTM update rule and use it for image captioning.

In [1]:

```
1  # As usual, a bit of setup
2  from __future__ import print_function
3  import time, os, json
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import nltk
7
8  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
9  from cs231n.rnn_layers import *
10 from cs231n.captioning_solver import CaptioningSolver
11 from cs231n.classifiers.rnn import CaptioningRNN
12 from cs231n.coco_utils import load_coco_data, sample_coco_minibatch, decode_captions
13 from cs231n.image_utils import image_from_url
14
15 %matplotlib inline
16 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
17 plt.rcParams['image.interpolation'] = 'nearest'
18 plt.rcParams['image.cmap'] = 'gray'
19
20 # for auto-reloading external modules
21 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-python
22 %load_ext autoreload
23 %autoreload 2
24
25 def rel_error(x, y):
26     """ returns relative error """
27     return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

Load MS-COCO data

As in the previous notebook, we will use the Microsoft COCO dataset for captioning.

In [2]:

```
1 # Load COCO data from disk; this returns a dictionary
2 # We'll work with dimensionality-reduced features for this notebook, but feel
3 # free to experiment with the original features by changing the flag below.
4 data = load_coco_data(pca_features=True)
5
6 # Print out all the keys and values from the data dictionary
7 for k, v in data.items():
8     if type(v) == np.ndarray:
9         print(k, type(v), v.shape, v.dtype)
10    else:
11        print(k, type(v), len(v))
```

```
train_captions <class 'numpy.ndarray'> (400135, 17) int32
train_image_idxes <class 'numpy.ndarray'> (400135,) int32
val_captions <class 'numpy.ndarray'> (195954, 17) int32
val_image_idxes <class 'numpy.ndarray'> (195954,) int32
train_features <class 'numpy.ndarray'> (82783, 512) float32
val_features <class 'numpy.ndarray'> (40504, 512) float32
idx_to_word <class 'list'> 1004
word_to_idx <class 'dict'> 1004
train_urls <class 'numpy.ndarray'> (82783,) <U63
val_urls <class 'numpy.ndarray'> (40504,) <U63
```

LSTM

If you read recent papers, you'll see that many people use a variant on the vanilla RNN called Long-Short Term Memory (LSTM) RNNs. Vanilla RNNs can be tough to train on long sequences due to vanishing and exploding gradients caused by repeated matrix multiplication. LSTMs solve this problem by replacing the simple update rule of the vanilla RNN with a gating mechanism as follows.

Similar to the vanilla RNN, at each timestep we receive an input $x_t \in \mathbb{R}^D$ and the previous hidden state $h_{t-1} \in \mathbb{R}^H$; the LSTM also maintains an H -dimensional *cell state*, so we also receive the previous cell state $c_{t-1} \in \mathbb{R}^H$. The learnable parameters of the LSTM are an *input-to-hidden* matrix $W_x \in \mathbb{R}^{4H \times D}$, a *hidden-to-hidden* matrix $W_h \in \mathbb{R}^{4H \times H}$ and a *bias vector* $b \in \mathbb{R}^{4H}$.

At each timestep we first compute an *activation vector* $a \in \mathbb{R}^{4H}$ as $a = W_x x_t + W_h h_{t-1} + b$. We then divide this into four vectors $a_i, a_f, a_o, a_g \in \mathbb{R}^H$ where a_i consists of the first H elements of a , a_f is the next H elements of a , etc. We then compute the *input gate* $g \in \mathbb{R}^H$, *forget gate* $f \in \mathbb{R}^H$, *output gate* $o \in \mathbb{R}^H$ and *block input* $g \in \mathbb{R}^H$ as

$$i = \sigma(a_i) \quad f = \sigma(a_f) \quad o = \sigma(a_o) \quad g = \tanh(a_g)$$

where σ is the sigmoid function and \tanh is the hyperbolic tangent, both applied elementwise.

Finally we compute the next cell state c_t and next hidden state h_t as

$$c_t = f \odot c_{t-1} + i \odot g \quad h_t = o \odot \tanh(c_t)$$

where \odot is the elementwise product of vectors.

In the rest of the notebook we will implement the LSTM update rule and apply it to the image captioning task.

In the code, we assume that data is stored in batches so that $X_t \in \mathbb{R}^{N \times D}$, and will work with *transposed* versions of the parameters: $W_x \in \mathbb{R}^{D \times 4H}$, $W_h \in \mathbb{R}^{H \times 4H}$ so that activations $A \in \mathbb{R}^{N \times 4H}$ can be computed efficiently as $A = X_t W_x + H_{t-1} W_h$

LSTM: step forward

Implement the forward pass for a single timestep of an LSTM in the `lstm_step_forward` function in the file `cs231n/rnn_layers.py`. This should be similar to the `rnn_step_forward` function that you implemented above, but using the LSTM update rule instead.

Once you are done, run the following to perform a simple test of your implementation. You should see errors around $1e-8$ or less.

In [3]:

```
1 N, D, H = 3, 4, 5
2 x = np.linspace(-0.4, 1.2, num=N*D).reshape(N, D)
3 prev_h = np.linspace(-0.3, 0.7, num=N*H).reshape(N, H)
4 prev_c = np.linspace(-0.4, 0.9, num=N*H).reshape(N, H)
5 Wx = np.linspace(-2.1, 1.3, num=4*D*H).reshape(D, 4 * H)
6 Wh = np.linspace(-0.7, 2.2, num=4*H*H).reshape(H, 4 * H)
7 b = np.linspace(0.3, 0.7, num=4*H)
8
9 next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)
10
11 expected_next_h = np.asarray([
12     [ 0.24635157,  0.28610883,  0.32240467,  0.35525807,  0.38474904],
13     [ 0.49223563,  0.55611431,  0.61507696,  0.66844003,  0.7159181 ],
14     [ 0.56735664,  0.66310127,  0.74419266,  0.80889665,  0.858299  ]])
15 expected_next_c = np.asarray([
16     [ 0.32986176,  0.39145139,  0.451556,    0.51014116,  0.56717407],
17     [ 0.66382255,  0.76674007,  0.87195994,  0.97902709,  1.08751345],
18     [ 0.74192008,  0.90592151,  1.07717006,  1.25120233,  1.42395676]])
19
20 print('next_h error: ', rel_error(expected_next_h, next_h))
21 print('next_c error: ', rel_error(expected_next_c, next_c))
```

next_h error: 5.7054131185818695e-09

next_c error: 5.8143123088804145e-09

LSTM: step backward

Implement the backward pass for a single LSTM timestep in the function `lstm_step_backward` in the file `cs231n/rnn_layers.py`. Once you are done, run the following to perform numeric gradient checking on your implementation. You should see errors around $1e-6$ or less.

In [4]:

```
1 np.random.seed(231)
2
3 N, D, H = 4, 5, 6
4 x = np.random.randn(N, D)
5 prev_h = np.random.randn(N, H)
6 prev_c = np.random.randn(N, H)
7 Wx = np.random.randn(D, 4 * H)
8 Wh = np.random.randn(H, 4 * H)
9 b = np.random.randn(4 * H)
10
11 next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)
12
13 dnext_h = np.random.randn(*next_h.shape)
14 dnext_c = np.random.randn(*next_c.shape)
15
16 fx_h = lambda x: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
17 fh_h = lambda h: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
18 fc_h = lambda c: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
19 fWx_h = lambda Wx: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
20 fWh_h = lambda Wh: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
21 fb_h = lambda b: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
22
23 fx_c = lambda x: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
24 fh_c = lambda h: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
25 fc_c = lambda c: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
26 fWx_c = lambda Wx: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
27 fWh_c = lambda Wh: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
28 fb_c = lambda b: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
29
30 num_grad = eval_numerical_gradient_array
31
32 dx_num = num_grad(fx_h, x, dnext_h) + num_grad(fx_c, x, dnext_c)
33 dh_num = num_grad(fh_h, prev_h, dnext_h) + num_grad(fh_c, prev_h, dnext_c)
34 dc_num = num_grad(fc_h, prev_c, dnext_h) + num_grad(fc_c, prev_c, dnext_c)
35 dWx_num = num_grad(fWx_h, Wx, dnext_h) + num_grad(fWx_c, Wx, dnext_c)
36 dWh_num = num_grad(fWh_h, Wh, dnext_h) + num_grad(fWh_c, Wh, dnext_c)
37 db_num = num_grad(fb_h, b, dnext_h) + num_grad(fb_c, b, dnext_c)
38
39 dx, dh, dc, dWx, dWh, db = lstm_step_backward(dnext_h, dnext_c, cache)
40
41 print('dx error: ', rel_error(dx_num, dx))
42 print('dh error: ', rel_error(dh_num, dh))
43 print('dc error: ', rel_error(dc_num, dc))
44 print('dWx error: ', rel_error(dWx_num, dWx))
45 print('dWh error: ', rel_error(dWh_num, dWh))
46 print('db error: ', rel_error(db_num, db))
```

```
dx error: 6.141307149471403e-10
dh error: 3.3953235055372503e-10
dc error: 1.5221771913099803e-10
dWx error: 1.6933643922734908e-09
dWh error: 2.5561308517943814e-08
db error: 1.7349247160222088e-10
```

LSTM: forward

In the function `lstm_forward` in the file `cs231n/rnn_layers.py`, implement the `lstm_forward` function to run an LSTM forward on an entire timeseries of data.

When you are done, run the following to check your implementation. You should see an error around $1e-7$.

In [5]:

```
1 N, D, H, T = 2, 5, 4, 3
2 x = np.linspace(-0.4, 0.6, num=N*T*D).reshape(N, T, D)
3 h0 = np.linspace(-0.4, 0.8, num=N*H).reshape(N, H)
4 Wx = np.linspace(-0.2, 0.9, num=4*D*H).reshape(D, 4 * H)
5 Wh = np.linspace(-0.3, 0.6, num=4*H*H).reshape(H, 4 * H)
6 b = np.linspace(0.2, 0.7, num=4*H)
7
8 h, cache = lstm_forward(x, h0, Wx, Wh, b)
9
10 expected_h = np.asarray([
11     [[ 0.01764008,  0.01823233,  0.01882671,  0.0194232 ],
12      [ 0.11287491,  0.12146228,  0.13018446,  0.13902939],
13      [ 0.31358768,  0.33338627,  0.35304453,  0.37250975]],
14     [[ 0.45767879,  0.4761092,   0.4936887,   0.51041945],
15      [ 0.6704845,   0.69350089,  0.71486014,  0.7346449 ],
16      [ 0.81733511,  0.83677871,  0.85403753,  0.86935314]]])
17
18 print('h error: ', rel_error(expected_h, h))
```

h error: 8.610537452106624e-08

LSTM: backward

Implement the backward pass for an LSTM over an entire timeseries of data in the function `lstm_backward` in the file `cs231n/rnn_layers.py`. When you are done, run the following to perform numeric gradient checking on your implementation. You should see errors around $1e-7$ or less.

In [6]:

```
1 from cs231n.rnn_layers import lstm_forward, lstm_backward
2 np.random.seed(231)
3
4 N, D, T, H = 2, 3, 10, 6
5
6 x = np.random.randn(N, T, D)
7 h0 = np.random.randn(N, H)
8 Wx = np.random.randn(D, 4 * H)
9 Wh = np.random.randn(H, 4 * H)
10 b = np.random.randn(4 * H)
11
12 out, cache = lstm_forward(x, h0, Wx, Wh, b)
13
14 dout = np.random.randn(*out.shape)
15
16 dx, dh0, dWx, dWh, db = lstm_backward(dout, cache)
17
18 fx = lambda x: lstm_forward(x, h0, Wx, Wh, b)[0]
19 fh0 = lambda h0: lstm_forward(x, h0, Wx, Wh, b)[0]
20 fWx = lambda Wx: lstm_forward(x, h0, Wx, Wh, b)[0]
21 fWh = lambda Wh: lstm_forward(x, h0, Wx, Wh, b)[0]
22 fb = lambda b: lstm_forward(x, h0, Wx, Wh, b)[0]
23
24 dx_num = eval_numerical_gradient_array(fx, x, dout)
25 dh0_num = eval_numerical_gradient_array(fh0, h0, dout)
26 dWx_num = eval_numerical_gradient_array(fWx, Wx, dout)
27 dWh_num = eval_numerical_gradient_array(fWh, Wh, dout)
28 db_num = eval_numerical_gradient_array(fb, b, dout)
29
30 print('dx error: ', rel_error(dx_num, dx))
31 print('dh0 error: ', rel_error(dh0_num, dh0))
32 print('dWx error: ', rel_error(dWx_num, dWx))
33 print('dWh error: ', rel_error(dWh_num, dWh))
34 print('db error: ', rel_error(db_num, db))
```

```
dx error: 4.825065794768194e-09
dh0 error: 7.500950523704672e-09
dWx error: 1.751994908422919e-09
dWh error: 1.0853770675498594e-06
db error: 7.42754365004995e-10
```

LSTM captioning model

Now that you have implemented an LSTM, update the implementation of the `loss` method of the `CaptioningRNN` class in the file `cs231n/classifiers/rnn.py` to handle the case where `self.cell_type` is `lstm`. This should require adding less than 10 lines of code.

Once you have done so, run the following to check your implementation. You should see a difference of less than $1e-10$.

In [7]:

```
1 N, D, W, H = 10, 20, 30, 40
2 word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
3 V = len(word_to_idx)
4 T = 13
5
6 model = CaptioningRNN(word_to_idx,
7                       input_dim=D,
8                       wordvec_dim=W,
9                       hidden_dim=H,
10                      cell_type='lstm',
11                      dtype=np.float64)
12
13 # Set all model parameters to fixed values
14 for k, v in model.params.items():
15     model.params[k] = np.linspace(-1.4, 1.3, num=v.size).reshape(*v.shape)
16
17 features = np.linspace(-0.5, 1.7, num=N*D).reshape(N, D)
18 captions = (np.arange(N * T) % V).reshape(N, T)
19
20 loss, grads = model.loss(features, captions)
21 expected_loss = 9.82445935443
22
23 print('loss: ', loss)
24 print('expected loss: ', expected_loss)
25 print('difference: ', abs(loss - expected_loss))
```

loss: 9.82445935443226

expected loss: 9.82445935443

difference: 2.261302256556519e-12

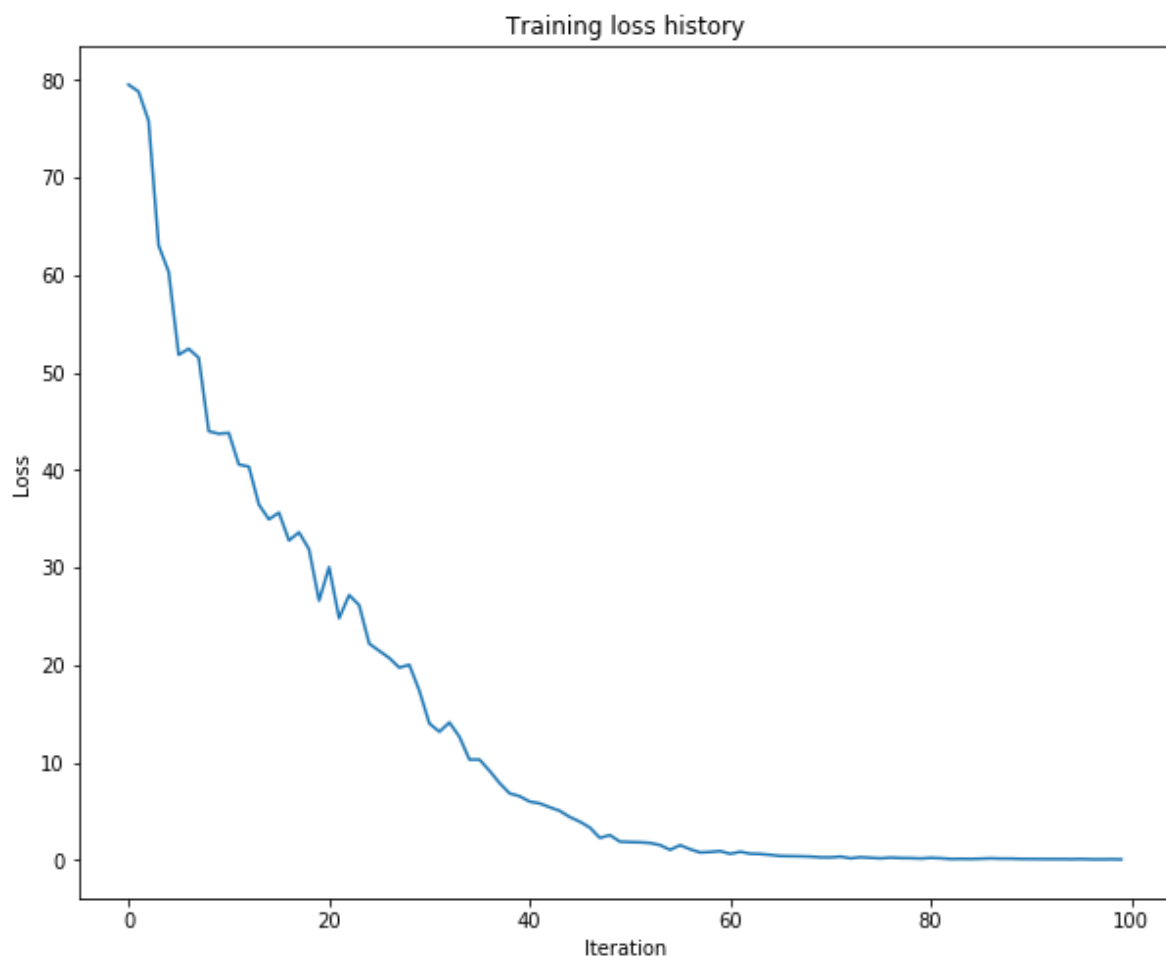
Overfit LSTM captioning model

Run the following to overfit an LSTM captioning model on the same small dataset as we used for the RNN previously. You should see losses less than 0.5.

In [8]:

```
1 np.random.seed(231)
2
3 small_data = load_coco_data(max_train=50)
4
5 small_lstm_model = CaptioningRNN(
6     cell_type='lstm',
7     word_to_idx=data['word_to_idx'],
8     input_dim=data['train_features'].shape[1],
9     hidden_dim=512,
10    wordvec_dim=256,
11    dtype=np.float32,
12)
13
14 small_lstm_solver = CaptioningSolver(small_lstm_model, small_data,
15     update_rule='adam',
16     num_epochs=50,
17     batch_size=25,
18     optim_config={
19         'learning_rate': 5e-3,
20     },
21     lr_decay=0.995,
22     verbose=True, print_every=10,
23)
24
25 small_lstm_solver.train()
26
27 # Plot the training losses
28 plt.plot(small_lstm_solver.loss_history)
29 plt.xlabel('Iteration')
30 plt.ylabel('Loss')
31 plt.title('Training loss history')
32 plt.show()
```

```
(Iteration 1 / 100) loss: 79.551152
(Iteration 11 / 100) loss: 43.829102
(Iteration 21 / 100) loss: 30.062537
(Iteration 31 / 100) loss: 14.020034
(Iteration 41 / 100) loss: 6.010117
(Iteration 51 / 100) loss: 1.856017
(Iteration 61 / 100) loss: 0.647575
(Iteration 71 / 100) loss: 0.286187
(Iteration 81 / 100) loss: 0.238507
(Iteration 91 / 100) loss: 0.122130
```

LSTM test-time sampling

Modify the `sample` method of the `CaptioningRNN` class to handle the case where `self.cell_type` is `lstm`. This should take fewer than 10 lines of code.

When you are done run the following to sample from your overfit LSTM model on some training and validation set samples.

In [9]:

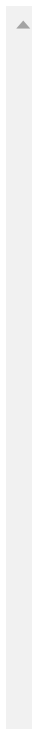
```
1 for split in ['train', 'val']:
2     minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2)
3     gt_captions, features, urls = minibatch
4     gt_captions = decode_captions(gt_captions, data['idx_to_word'])
5
6     sample_captions = small_lstm_model.sample(features)
7     sample_captions = decode_captions(sample_captions, data['idx_to_word'])
8
9     for gt_caption, sample_caption, url in zip(gt_captions, sample_captions, urls):
10         plt.imshow(image_from_url(url))
11         plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
12         plt.axis('off')
13         plt.show()
```

train
<START> a man standing on the side of a road with bags of luggage <END>
GT:<START> a man standing on the side of a road with bags of luggage <END>



train

<START> a man <UNK> with a bright colorful kite <END>
GT:<START> a man <UNK> with a bright colorful kite <END>



val
<START> a person <UNK> of a <UNK> <END>
GT:<START> a sign that is on the front of a train station <END>



val
<START> a cat is <UNK> and a <UNK> <END>
GT:<START> a car is parked on a street at night <END>



Train a good captioning model (extra credit for 4803)

Using the pieces you have implemented in this and the previous notebook, train a captioning model that gives decent qualitative results (better than the random garbage you saw with the overfit models) when sampling on the validation set. You can subsample the training set if you want; we just want to see samples on the validation set that are better than random.

In addition to qualitatively evaluating your model by inspecting its results, you can also quantitatively evaluate your model using the BLEU unigram precision metric. In order to achieve full credit you should train a model that achieves a BLEU unigram score of >0.3 . BLEU scores range from 0 to 1; the closer to 1, the better. Here's a reference to the [paper](http://www.aclweb.org/anthology/P02-1040.pdf) (<http://www.aclweb.org/anthology/P02-1040.pdf>) that introduces BLEU if you're interested in learning more about how it works.

Feel free to use PyTorch for this section if you'd like to train faster on a GPU... though you can definitely get above 0.3 using your Numpy code. We're providing you the evaluation code that is compatible with the Numpy model as defined above... you should be able to adapt it for PyTorch if you go that route.

Create the model in the file `cs231n/classifiers/mymodel.py`. You can base it after the `CaptioningRNN` class. Write a text comment in the delineated cell below explaining what you tried in your model.

Also add a cell below that trains and tests your model. Make sure to include the call to `evaluate_model` which prints out your highest validation BLEU score for full credit.

In [11]:

```
1 def BLEU_score(gt_caption, sample_caption):
2     """
3     gt_caption: string, ground-truth caption
4     sample_caption: string, your model's predicted caption
5     Returns unigram BLEU score.
6     """
7     reference = [x for x in gt_caption.split(' ')
8                   if ('<END>' not in x and '<START>' not in x and '<UNK>' not in x)]
9     hypothesis = [x for x in sample_caption.split(' ')
10                  if ('<END>' not in x and '<START>' not in x and '<UNK>' not in x)]
11     BLEUScore = nltk.translate.bleu_score.sentence_bleu([reference], hypothesis, weights = [1])
12     return BLEUScore
13
14 def evaluate_model(model):
15     """
16     model: CaptioningRNN model
17     Prints unigram BLEU score averaged over 1000 training and val examples.
18     """
19     BLEUScores = {}
20     for split in ['train', 'val']:
21         minibatch = sample_coco_minibatch(data, split=split, batch_size=1000)
22         gt_captions, features, urls = minibatch
23         gt_captions = decode_captions(gt_captions, data['idx_to_word'])
24
25         sample_captions = model.sample(features)
26         sample_captions = decode_captions(sample_captions, data['idx_to_word'])
27
28         total_score = 0.0
29         for gt_caption, sample_caption, url in zip(gt_captions, sample_captions, urls):
30             total_score += BLEU_score(gt_caption, sample_caption)
31
32         BLEUScores[split] = total_score / len(sample_captions)
33
34     for split in BLEUScores:
35         print('Average BLEU score for %s: %f' % (split, BLEUScores[split]))
```

write a description of your model here:

First, the dataset I choose is 25000 based on my computational resource. In order to increase the training speed, I choose 2048 and 1024 as the batch size. Learning rate is very small. I tried many different number and finally choose $5e-3$.

In [17]:

```
1  # write your code to train your model here.
2  # make sure to include the call to evaluate_model which prints out your highest validation BLEU
3  from cs231n.classifiers.mymodel import Mymodel
4
5  np.random.seed(231)
6
7  small_data = load_coco_data(max_train=25000)
8
9  small_lstm_model = Mymodel(
10     cell_type='lstm',
11     word_to_idx=data['word_to_idx'],
12     input_dim=data['train_features'].shape[1],
13     hidden_dim=1024,
14     wordvec_dim=2048,
15     dtype=np.float64,
16 )
17
18  small_lstm_solver = CaptioningSolver(small_lstm_model, small_data,
19     update_rule='adam',
20     num_epochs=20,
21     batch_size=32,
22     optim_config={
23         'learning_rate': 5e-3,
24     },
25     lr_decay=0.996,
26     verbose=True, print_every=1000,
27 )
28
29  small_lstm_solver.train()
30  evaluate_model(small_lstm_model)
```

```
(Iteration 1 / 15620) loss: 75.182688
(Iteration 1001 / 15620) loss: 25.802610
(Iteration 2001 / 15620) loss: 23.818434
(Iteration 3001 / 15620) loss: 19.109261
(Iteration 4001 / 15620) loss: 15.545158
(Iteration 5001 / 15620) loss: 14.839553
(Iteration 6001 / 15620) loss: 12.614585
(Iteration 7001 / 15620) loss: 9.706871
(Iteration 8001 / 15620) loss: 9.223613
(Iteration 9001 / 15620) loss: 9.511062
(Iteration 10001 / 15620) loss: 5.808891
(Iteration 11001 / 15620) loss: 5.041856
(Iteration 12001 / 15620) loss: 3.246330
(Iteration 13001 / 15620) loss: 4.050406
(Iteration 14001 / 15620) loss: 3.231652
(Iteration 15001 / 15620) loss: 2.824608
Average BLEU score for train: 0.309956
Average BLEU score for val: 0.304651
```