# Image Captioning with RNNs

In this exercise you will implement a vanilla recurrent neural networks and use them it to train a model that can generate novel captions for images.

## Install h5py

The COCO dataset we will be using is stored in HDF5 format. To load HDF5 files, we will need to install the `h5py` Python package. From the command line, run:

`pip install h5py`

If you receive a permissions error, you may need to run the command as root:

`sudo pip install h5py`

You can also run commands directly from the Jupyter notebook by prefixing the command with the "!" character:

In [2]:

```
!pip install h5py
```

```
Requirement already satisfied: h5py in /opt/anaconda3/lib/python3.7/site-packages (2.8.0)
Requirement already satisfied: numpy>=1.7 in /opt/anaconda3/lib/python3.7/site-packages (from
h5py) (1.15.4)
Requirement already satisfied: six in /opt/anaconda3/lib/python3.7/site-packages (from h5py)
(1.12.0)
```

## Look at the data

It is always a good idea to look at examples from the dataset before working with it.

You can use the `sample_coco_minibatch` function from the file `cs231n/coco_utils.py` to sample minibatches of data from the data structure returned from `load_coco_data`. Run the following to sample a small minibatch of training data and show the images and their captions. Running it multiple times and looking at the results helps you to get a sense of the dataset.

Note that we decode the captions using the `decode_captions` function and that we download the images on-the-fly using their Flickr URL, so **you must be connected to the internet to view images**.

In [4]:

```
# Sample a minibatch and show the images and captions
batch_size = 3

captions, features, urls = sample_coco_minibatch(data, batch_size=batch_size)
for i, (caption, url) in enumerate(zip(captions, urls)):
    plt.imshow(image_from_url(url))
    plt.axis('off')
    caption_str = decode_captions(caption, data['idx_to_word'])
    plt.title(caption_str)
    plt.show()
```

<START> a baseball player holding a bat standing next to home plate <END>



<START> a room with wooden floors and red walls <END>

<START> the <UNK> is <UNK> under the microwave oven <END>



## Recurrent Neural Networks

As discussed in lecture, we will use recurrent neural network (RNN) language models for image captioning. The file `cs231n/rnn_layers.py` contains implementations of different layer types that are needed for recurrent neural networks, and the file `cs231n/classifiers/rnn.py` uses these layers to implement an image captioning model.

We will first implement different types of RNN layers in `cs231n/rnn_layers.py`.

## Vanilla RNN: step forward

Open the file `cs231n/rnn_layers.py`. This file implements the forward and backward passes for different types of layers that are commonly used in recurrent neural networks.

First implement the function `rnn_step_forward` which implements the forward pass for a single timestep of a vanilla recurrent neural network. After doing so run the following to check your implementation. You should see errors on the order of e-8 or less.

In [5]:

```
N, D, H = 3, 10, 4

x = np.linspace(-0.4, 0.7, num=N*D).reshape(N, D)
prev_h = np.linspace(-0.2, 0.5, num=N*H).reshape(N, H)
Wx = np.linspace(-0.1, 0.9, num=D*H).reshape(D, H)
Wh = np.linspace(-0.3, 0.7, num=H*H).reshape(H, H)
b = np.linspace(-0.2, 0.4, num=H)

next_h, _ = rnn_step_forward(x, prev_h, Wx, Wh, b)
expected_next_h = np.asarray([
  [-0.58172089, -0.50182032, -0.41232771, -0.31410098],
  [ 0.66854692,  0.79562378,  0.87755553,  0.92795967],
  [ 0.97934501,  0.99144213,  0.99646691,  0.99854353]])

print('next_h error: ', rel_error(expected_next_h, next_h))
```

```
next_h error:  6.292421426471037e-09
```

## Vanilla RNN: step backward

In the file `cs231n/rnn_layers.py` implement the `rnn_step_backward` function. After doing so run the following to numerically gradient check your implementation. You should see errors on the order of `e-8` or less.

In [6]:

```python
from cs231n.rnn_layers import rnn_step_forward, rnn_step_backward
np.random.seed(231)
N, D, H = 4, 5, 6
x = np.random.randn(N, D)
h = np.random.randn(N, H)
Wx = np.random.randn(D, H)
Wh = np.random.randn(H, H)
b = np.random.randn(H)

out, cache = rnn_step_forward(x, h, Wx, Wh, b)

dnext_h = np.random.randn(*out.shape)

fx = lambda x: rnn_step_forward(x, h, Wx, Wh, b)[0]
fh = lambda prev_h: rnn_step_forward(x, h, Wx, Wh, b)[0]
fWx = lambda Wx: rnn_step_forward(x, h, Wx, Wh, b)[0]
fWh = lambda Wh: rnn_step_forward(x, h, Wx, Wh, b)[0]
fb = lambda b: rnn_step_forward(x, h, Wx, Wh, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dnext_h)
dprev_h_num = eval_numerical_gradient_array(fh, h, dnext_h)
dWx_num = eval_numerical_gradient_array(fWx, Wx, dnext_h)
dWh_num = eval_numerical_gradient_array(fWh, Wh, dnext_h)
db_num = eval_numerical_gradient_array(fb, b, dnext_h)

dx, dprev_h, dWx, dWh, db = rnn_step_backward(dnext_h, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dprev_h error: ', rel_error(dprev_h_num, dprev_h))
print('dWx error: ', rel_error(dWx_num, dWx))
print('dWh error: ', rel_error(dWh_num, dWh))
print('db error: ', rel_error(db_num, db))
```

```
dx error:  4.0192769090159184e-10
dprev_h error:  2.5136656668664053e-10
dWx error:  3.398875305713782e-10
dWh error:  3.355162782632426e-10
db error:  1.946925061042176e-10
```

# Vanilla RNN: forward

Now that you have implemented the forward and backward passes for a single timestep of a vanilla RNN, you will combine these pieces to implement a RNN that processes an entire sequence of data.

In the file `cs231n/rnn_layers.py`, implement the function `rnn_forward`. This should be implemented using the `rnn_step_forward` function that you defined above. After doing so run the following to check your implementation. You should see errors on the order of `e-7` or less.

In [7]:

```python
N, T, D, H = 2, 3, 4, 5

x = np.linspace(-0.1, 0.3, num=N*T*D).reshape(N, T, D)
h0 = np.linspace(-0.3, 0.1, num=N*H).reshape(N, H)
Wx = np.linspace(-0.2, 0.4, num=D*H).reshape(D, H)
Wh = np.linspace(-0.4, 0.1, num=H*H).reshape(H, H)
b = np.linspace(-0.7, 0.1, num=H)

h, _ = rnn_forward(x, h0, Wx, Wh, b)
expected_h = np.asarray([
  [
    [-0.42070749, -0.27279261, -0.11074945,  0.05740409,  0.22236251],
    [-0.39525808, -0.22554661, -0.0409454,   0.14649412,  0.32397316],
    [-0.42305111, -0.24223728, -0.04287027,  0.15997045,  0.35014525],
  ],
  [
    [-0.55857474, -0.39065825, -0.19198182,  0.02378408,  0.23735671],
    [-0.27150199, -0.07088804,  0.13562939,  0.33099728,  0.50158768],
    [-0.51014825, -0.30524429, -0.06755202,  0.17806392,  0.40333043]])
```

```
                   [ 0.31011023,  0.30321129,  0.00733202,  0.17000392,  0.18333313]]])
print('h error: ', rel_error(expected_h, h))
```

```
h error:  7.728466180186066e-08
```

## Vanilla RNN: backward

In the file `cs231n/rnn_layers.py`, implement the backward pass for a vanilla RNN in the function `rnn_backward`. This should run back-propagation over the entire sequence, making calls to the `rnn_step_backward` function that you defined earlier. You should see errors on the order of e-6 or less.

In [8]:

```
np.random.seed(231)

N, D, T, H = 2, 3, 10, 5

x = np.random.randn(N, T, D)
h0 = np.random.randn(N, H)
Wx = np.random.randn(D, H)
Wh = np.random.randn(H, H)
b = np.random.randn(H)

out, cache = rnn_forward(x, h0, Wx, Wh, b)

dout = np.random.randn(*out.shape)

dx, dh0, dWx, dWh, db = rnn_backward(dout, cache)

fx = lambda x: rnn_forward(x, h0, Wx, Wh, b)[0]
fh0 = lambda h0: rnn_forward(x, h0, Wx, Wh, b)[0]
fWx = lambda Wx: rnn_forward(x, h0, Wx, Wh, b)[0]
fWh = lambda Wh: rnn_forward(x, h0, Wx, Wh, b)[0]
fb = lambda b: rnn_forward(x, h0, Wx, Wh, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
dh0_num = eval_numerical_gradient_array(fh0, h0, dout)
dWx_num = eval_numerical_gradient_array(fWx, Wx, dout)
dWh_num = eval_numerical_gradient_array(fWh, Wh, dout)
db_num = eval_numerical_gradient_array(fb, b, dout)

print('dx error: ', rel_error(dx_num, dx))
print('dh0 error: ', rel_error(dh0_num, dh0))
print('dWx error: ', rel_error(dWx_num, dWx))
print('dWh error: ', rel_error(dWh_num, dWh))
print('db error: ', rel_error(db_num, db))
```

```
dx error:  2.12371994872961e-09
dh0 error:  3.380520197084487e-09
dWx error:  7.133880725895019e-09
dWh error:  1.2991706887909817e-07
db error:  4.309473374164083e-10
```

## Word embedding: forward

In deep learning systems, we commonly represent words using vectors. Each word of the vocabulary will be associated with a vector, and these vectors will be learned jointly with the rest of the system.

In the file `cs231n/rnn_layers.py`, implement the function `word_embedding_forward` to convert words (represented by integers) into vectors. Run the following to check your implementation. You should see an error on the order of `e-8` or less.

In [9]:

```
N, T, V, D = 2, 4, 5, 3

x = np.asarray([[0, 3, 1, 2], [2, 1, 0, 3]])
W = np.linspace(0, 1, num=V*D).reshape(V, D)

out, _ = word_embedding_forward(x, W)
expected_out = np.asarray([
 [[ 0.,          0.07142857,  0.14285714],
  [ 0.64285714,  0.71428571,  0.78571429],
  [ 0.21428571,  0.28571429,  0.35714286],
```

```
   [ 0.42857143,  0.5,         0.57142857]],
 [[ 0.42857143,  0.5,         0.57142857],
  [ 0.21428571,  0.28571429,  0.35714286],
  [ 0.,          0.07142857,  0.14285714],
  [ 0.64285714,  0.71428571,  0.78571429]]])

print('out error: ', rel_error(expected_out, out))
```

```
out error:  1.000000094736443e-08
```

## Word embedding: backward

Implement the backward pass for the word embedding function in the function `word_embedding_backward`. After doing so run the following to numerically gradient check your implementation. You should see an error on the order of `e-11` or less.

```
np.random.seed(231)

N, T, V, D = 50, 3, 5, 6
x = np.random.randint(V, size=(N, T))
W = np.random.randn(V, D)

out, cache = word_embedding_forward(x, W)
dout = np.random.randn(*out.shape)
dW = word_embedding_backward(dout, cache)

f = lambda W: word_embedding_forward(x, W)[0]
dW_num = eval_numerical_gradient_array(f, W, dout)

print('dW error: ', rel_error(dW, dW_num))
```

```
dW error:  3.2774595693100364e-12
```

## RNN for image captioning

Now that you have implemented the necessary layers, you can combine them to build an image captioning model. Open the file `cs231n/classifiers/rnn.py` and look at the `CaptioningRNN` class.

Implement the forward and backward pass of the model in the `loss` function. For now you only need to implement the case where `cell_type='rnn'` for vanialla RNNs; you will implement the LSTM case later. After doing so, run the following to check your forward pass using a small test case; you should see error on the order of `e-10` or less.

```
N, D, W, H = 10, 20, 30, 40
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
V = len(word_to_idx)
T = 13

model = CaptioningRNN(word_to_idx,
          input_dim=D,
          wordvec_dim=W,
          hidden_dim=H,
          cell_type='rnn',
          dtype=np.float64)

# Set all model parameters to fixed values
for k, v in model.params.items():
    model.params[k] = np.linspace(-1.4, 1.3, num=v.size).reshape(*v.shape)

features = np.linspace(-1.5, 0.3, num=(N * D)).reshape(N, D)
captions = (np.arange(N * T) % V).reshape(N, T)

loss, grads = model.loss(features, captions)
expected_loss = 9.83235591003

print('loss: ', loss)
print('expected loss: ', expected_loss)
print('difference: ', abs(loss - expected_loss))
```

```
loss:  9.832355910027388
```

```
        loss:   9.88238591884788
expected loss:   9.83235591003
difference:   2.611244553918368e-12
```

Run the following cell to perform numeric gradient checking on the `CaptioningRNN` class; you should see errors around the order of `e-6` or less.

```
np.random.seed(231)

batch_size = 2
timesteps = 3
input_dim = 4
wordvec_dim = 5
hidden_dim = 6
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
vocab_size = len(word_to_idx)

captions = np.random.randint(vocab_size, size=(batch_size, timesteps))
features = np.random.randn(batch_size, input_dim)

model = CaptioningRNN(word_to_idx,
          input_dim=input_dim,
          wordvec_dim=wordvec_dim,
          hidden_dim=hidden_dim,
          cell_type='rnn',
          dtype=np.float64,
        )

loss, grads = model.loss(features, captions)

for param_name in sorted(grads):
    f = lambda _: model.loss(features, captions)[0]
    param_grad_num = eval_numerical_gradient(f, model.params[param_name], verbose=False, h=1e-6)
    e = rel_error(param_grad_num, grads[param_name])
    print('%s relative error: %e' % (param_name, e))
```

```
W_embed relative error: 2.331074e-09
W_proj relative error: 9.974427e-09
W_vocab relative error: 2.875061e-09
Wh relative error: 4.685196e-09
Wx relative error: 7.725620e-07
b relative error: 4.909225e-10
b_proj relative error: 1.934808e-08
b_vocab relative error: 1.781169e-09
```

# Overfit small data

Similar to the `Solver` class that we used to train image classification models on the previous assignment, on this assignment we use a `CaptioningSolver` class to train image captioning models. Open the file `cs231n/captioning_solver.py` and read through the `CaptioningSolver` class; it should look very familiar.

Once you have familiarized yourself with the API, run the following to make sure your model overfits a small sample of 100 training examples. You should see a final loss of less than 0.1.

```
np.random.seed(231)

small_data = load_coco_data(max_train=50)

small_rnn_model = CaptioningRNN(
          cell_type='rnn',
          word_to_idx=data['word_to_idx'],
          input_dim=data['train_features'].shape[1],
          hidden_dim=512,
          wordvec_dim=256,
        )

small_rnn_solver = CaptioningSolver(small_rnn_model, small_data,
          update_rule='adam',
          num_epochs=50,
```

```
            batcn_size=25,
            optim_config={
                'learning_rate': 5e-3,
            },
            lr_decay=0.95,
            verbose=True, print_every=10,
        )

small_rnn_solver.train()

# Plot the training losses
plt.plot(small_rnn_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```
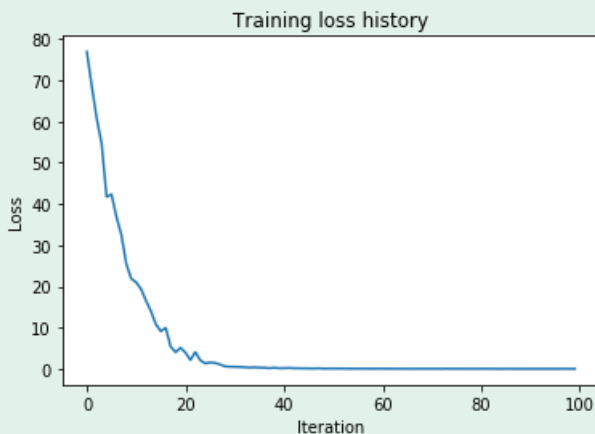
```
(Iteration 1 / 100) loss: 76.913487
(Iteration 11 / 100) loss: 21.063202
(Iteration 21 / 100) loss: 4.016187
(Iteration 31 / 100) loss: 0.567069
(Iteration 41 / 100) loss: 0.239435
(Iteration 51 / 100) loss: 0.162025
(Iteration 61 / 100) loss: 0.111542
(Iteration 71 / 100) loss: 0.097584
(Iteration 81 / 100) loss: 0.099099
(Iteration 91 / 100) loss: 0.073980
```



# Test-time sampling

Unlike classification models, image captioning models behave very differently at training time and at test time. At training time, we have access to the ground-truth caption, so we feed ground-truth words as input to the RNN at each timestep. At test time, we sample from the distribution over the vocabulary at each timestep, and feed the sample as input to the RNN at the next timestep.

In the file `cs231n/classifiers/rnn.py`, implement the `sample` method for test-time sampling. After doing so, run the following to sample from your overfitted model on both training and validation data. The samples on training data should be very good; the samples on validation data probably won't make sense.

In [22]:

```
for split in ['train', 'val']:
    minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2)
    gt_captions, features, urls = minibatch
    gt_captions = decode_captions(gt_captions, data['idx_to_word'])

    sample_captions = small_rnn_model.sample(features)
    sample_captions = decode_captions(sample_captions, data['idx_to_word'])

    for gt_caption, sample_caption, url in zip(gt_captions, sample_captions, urls):
        plt.imshow(image_from_url(url))
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
        plt.axis('off')
        plt.show()
```

train
a plane flying close to the ground as <UNK> coming in for landing <END>
GT:<START> a plane flying close to the ground as <UNK> coming in for landing <END>

train
a boy sitting with <UNK> on with a donut in his hand <END>
GT:<START> a boy sitting with <UNK> on with a donut in his hand <END>



val
bags the <UNK> the <UNK> while in <END>
GT:<START> a picture of a giraffes head eating leaves off a tree <END>



val
glasses seat in colorful car <END>
GT:<START> two brown and white cows standing next to each other <END>



# INLINE QUESTION 1

In our current image captioning setup, our RNN language model produces a word at every timestep as its output. However, an alternate way to pose the problem is to train the network to operate over *characters* (e.g. 'a', 'b', etc.) as opposed to words, so that at it every timestep, it receives the previous character as input and tries to predict the next character in the sequence. For example, the network might generate a caption like

'A', ' ', 'c', 'a', 't', ' ', 'o', 'n', ' ', 'a', ' ', 'b', 'e', 'd'

Can you describe one advantage of an image-captioning model that uses a character-level RNN? Can you also describe one disadvantage? HINT: there are several valid answers, but it might be useful to compare the parameter space of word-level and character-level models.

**Your Answer:** The advantage is that the discrete space of character-level RNN is much smaller (the amount of characters is much less than the amount of words). The main disadvantage is the longer training time because the sequence length is increasing.

In [ ]:

# 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.

## Load MS-COCO data

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

In [2]:

```python
# Load COCO data from disk; this returns a dictionary
# We'll work with dimensionality-reduced features for this notebook, but feel
# free to experiment with the original features by changing the flag below.
data = load_coco_data(pca_features=True)

# Print out all the keys and values from the data dictionary
for k, v in data.items():
    if type(v) == np.ndarray:
        print(k, type(v), v.shape, v.dtype)
    else:
        print(k, type(v), len(v))
```

```
train_captions <class 'numpy.ndarray'> (400135, 17) int32
train_image_idxs <class 'numpy.ndarray'> (400135,) int32
val_captions <class 'numpy.ndarray'> (195954, 17) int32
val_image_idxs <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 $x_t \in \mathrm{R}^D$ and the previous hidden state $h_{t-1} \in \mathrm{R}^H$; the LSTM also maintains an $H$-dimensional *cell state*, so we also receive the previous cell state $c_{t-1} \in \mathrm{R}^H$. The learnable parameters of the LSTM are an *input-to-hidden* matrix $W_x \in \mathrm{R}^{4H \times D}$, a *hidden-to-hidden* matrix $W_h \in \mathrm{R}^{4H \times H}$ and a *bias vector* $b \in \mathrm{R}^{4H}$.

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

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

where $\sigma$ 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

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

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

## LSTM: step forward

# 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 on the order of `e-8` or less.

In [4]:

```
N, D, H = 3, 4, 5
x = np.linspace(-0.4, 1.2, num=N*D).reshape(N, D)
prev_h = np.linspace(-0.3, 0.7, num=N*H).reshape(N, H)
prev_c = np.linspace(-0.4, 0.9, num=N*H).reshape(N, H)
Wx = np.linspace(-2.1, 1.3, num=4*D*H).reshape(D, 4 * H)
Wh = np.linspace(-0.7, 2.2, num=4*H*H).reshape(H, 4 * H)
b = np.linspace(0.3, 0.7, num=4*H)

next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)

expected_next_h = np.asarray([
    [ 0.24635157,  0.28610883,  0.32240467,  0.35525807,  0.38474904],
    [ 0.49223563,  0.55611431,  0.61507696,  0.66844003,  0.7159181 ],
    [ 0.56735664,  0.66310127,  0.74419266,  0.80889665,  0.858299  ]])
expected_next_c = np.asarray([
    [ 0.32986176,  0.39145139,  0.451556,    0.51014116,  0.56717407],
    [ 0.66382255,  0.76674007,  0.87195994,  0.97902709,  1.08751345],
    [ 0.74192008,  0.90592151,  1.07717006,  1.25120233,  1.42395676]])

print('next_h error: ', rel_error(expected_next_h, next_h))
print('next_c error: ', rel_error(expected_next_c, next_c))
```

```
next_h error:  5.7054131967097955e-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 on the order of `e-7` or less.

In [5]:

```
np.random.seed(231)

N, D, H = 4, 5, 6
x = np.random.randn(N, D)
prev_h = np.random.randn(N, H)
prev_c = np.random.randn(N, H)
Wx = np.random.randn(D, 4 * H)
Wh = np.random.randn(H, 4 * H)
b = np.random.randn(4 * H)

next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)

dnext_h = np.random.randn(*next_h.shape)
dnext_c = np.random.randn(*next_c.shape)

fx_h = lambda x: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
fh_h = lambda h: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
fc_h = lambda c: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
fWx_h = lambda Wx: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
fWh_h = lambda Wh: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]
fb_h = lambda b: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[0]

fx_c = lambda x: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
fh_c = lambda h: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
fc_c = lambda c: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
fWx_c = lambda Wx: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
fWh_c = lambda Wh: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
fb_c = lambda b: lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b)[1]
```

```
num_grad = eval_numerical_gradient_array

dx_num = num_grad(fx_h, x, dnext_h) + num_grad(fx_c, x, dnext_c)
dh_num = num_grad(fh_h, prev_h, dnext_h) + num_grad(fh_c, prev_h, dnext_c)
dc_num = num_grad(fc_h, prev_c, dnext_h) + num_grad(fc_c, prev_c, dnext_c)
dWx_num = num_grad(fWx_h, Wx, dnext_h) + num_grad(fWx_c, Wx, dnext_c)
dWh_num = num_grad(fWh_h, Wh, dnext_h) + num_grad(fWh_c, Wh, dnext_c)
db_num = num_grad(fb_h, b, dnext_h) + num_grad(fb_c, b, dnext_c)

dx, dh, dc, dWx, dWh, db = lstm_step_backward(dnext_h, dnext_c, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dh error: ', rel_error(dh_num, dh))
print('dc error: ', rel_error(dc_num, dc))
print('dWx error: ', rel_error(dWx_num, dWx))
print('dWh error: ', rel_error(dWh_num, dWh))
print('db error: ', rel_error(db_num, db))
```

```
dx error:  6.335032254429549e-10
dh error:  3.3963774090592634e-10
dc error:  1.5221771913099803e-10
dWx error:  2.1010960934639614e-09
dWh error:  9.712296180612259e-08
db error:  2.4915214652298706e-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 on the order of `e-7` or less.

In [9]:

```
N, D, H, T = 2, 5, 4, 3
x = np.linspace(-0.4, 0.6, num=N*T*D).reshape(N, T, D)
h0 = np.linspace(-0.4, 0.8, num=N*H).reshape(N, H)
Wx = np.linspace(-0.2, 0.9, num=4*D*H).reshape(D, 4 * H)
Wh = np.linspace(-0.3, 0.6, num=4*H*H).reshape(H, 4 * H)
b = np.linspace(0.2, 0.7, num=4*H)

h, cache = lstm_forward(x, h0, Wx, Wh, b)

expected_h = np.asarray([
 [[ 0.01764008,  0.01823233,  0.01882671,  0.0194232 ],
  [ 0.11287491,  0.12146228,  0.13018446,  0.13902939],
  [ 0.31358768,  0.33338627,  0.35304453,  0.37250975]],
 [[ 0.45767879,  0.4761092,   0.4936887,   0.51041945],
  [ 0.6704845,   0.69350089,  0.71486014,  0.7346449 ],
  [ 0.81733511,  0.83677871,  0.85403753,  0.86935314]]])

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 on the order of `e-8` or less. (For `dWh`, it's fine if your error is on the order of `e-6` or less).

In [11]:

```
from cs231n.rnn_layers import lstm_forward, lstm_backward
np.random.seed(231)

N, D, T, H = 2, 3, 10, 6

x = np.random.randn(N, T, D)
h0 = np.random.randn(N, H)
Wx = np.random.randn(D, 4 * H)
Wh = np.random.randn(H, 4 * H)
```

```
b = np.random.randn(4 * H)

out, cache = lstm_forward(x, h0, Wx, Wh, b)

dout = np.random.randn(*out.shape)

dx, dh0, dWx, dWh, db = lstm_backward(dout, cache)

fx  = lambda x: lstm_forward(x, h0, Wx, Wh, b)[0]
fh0 = lambda h0: lstm_forward(x, h0, Wx, Wh, b)[0]
fWx = lambda Wx: lstm_forward(x, h0, Wx, Wh, b)[0]
fWh = lambda Wh: lstm_forward(x, h0, Wx, Wh, b)[0]
fb  = lambda b: lstm_forward(x, h0, Wx, Wh, b)[0]

dx_num  = eval_numerical_gradient_array(fx, x, dout)
dh0_num = eval_numerical_gradient_array(fh0, h0, dout)
dWx_num = eval_numerical_gradient_array(fWx, Wx, dout)
dWh_num = eval_numerical_gradient_array(fWh, Wh, dout)
db_num  = eval_numerical_gradient_array(fb, b, dout)

print('dx error: ', rel_error(dx_num, dx))
print('dh0 error: ', rel_error(dh0_num, dh0))
print('dWx error: ', rel_error(dWx_num, dWx))
print('dWh error: ', rel_error(dWh_num, dWh))
print('db error: ', rel_error(db_num, db))
```

```
dx error:   7.838507661490775e-09
dh0 error:  2.469093442332882e-08
dWx error:  4.748335917394258e-09
dWh error:  1.042440847968619e-06
db error:   1.915271675612951e-09
```

## INLINE QUESTION

Recall that in an LSTM the input gate $i$, forget gate $f$, and output gate $o$ are all outputs of a sigmoid function. Why don't we use the ReLU activation function instead of sigmoid to compute these values? Explain.

**Your Answer:** Because if we use relu, all the outputs from the cell, as well as the cell state, will be strictly >= 0. In this way gradients become extremely large and are exploding. Using sigmoid can limit the output to be in range [0, 1].

## 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 on the order of `e-10` or less.

In [12]:

```
N, D, W, H = 10, 20, 30, 40
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
V = len(word_to_idx)
T = 13

model = CaptioningRNN(word_to_idx,
          input_dim=D,
          wordvec_dim=W,
          hidden_dim=H,
          cell_type='lstm',
          dtype=np.float64)

# Set all model parameters to fixed values
for k, v in model.params.items():
  model.params[k] = np.linspace(-1.4, 1.3, num=v.size).reshape(*v.shape)

features = np.linspace(-0.5, 1.7, num=N*D).reshape(N, D)
captions = (np.arange(N * T) % V).reshape(N, T)

loss, grads = model.loss(features, captions)
```

```
expected_loss = 9.82445935443

print('loss: ', loss)
print('expected loss: ', expected_loss)
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 a final loss less than 0.5.

In [13]:

```
np.random.seed(231)

small_data = load_coco_data(max_train=50)

small_lstm_model = CaptioningRNN(
          cell_type='lstm',
          word_to_idx=data['word_to_idx'],
          input_dim=data['train_features'].shape[1],
          hidden_dim=512,
          wordvec_dim=256,
          dtype=np.float32,
        )

small_lstm_solver = CaptioningSolver(small_lstm_model, small_data,
          update_rule='adam',
          num_epochs=50,
          batch_size=25,
          optim_config={
             'learning_rate': 5e-3,
          },
          lr_decay=0.995,
          verbose=True, print_every=10,
        )

small_lstm_solver.train()

# Plot the training losses
plt.plot(small_lstm_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```
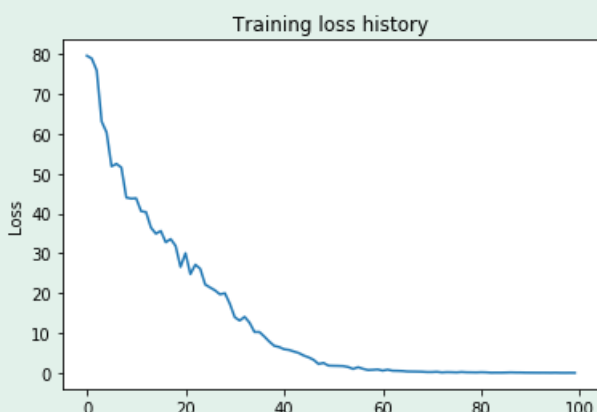
```
(Iteration 1 / 100) loss: 79.551150
(Iteration 11 / 100) loss: 43.829102
(Iteration 21 / 100) loss: 30.062640
(Iteration 31 / 100) loss: 14.020040
(Iteration 41 / 100) loss: 6.003100
(Iteration 51 / 100) loss: 1.853265
(Iteration 61 / 100) loss: 0.641483
(Iteration 71 / 100) loss: 0.284055
(Iteration 81 / 100) loss: 0.238315
(Iteration 91 / 100) loss: 0.124627
```

# 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. As with the RNN, training results should be very good, and validation results probably won't make a lot of sense (because we're overfitting).

In [15]:

```
for split in ['train', 'val']:
    minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2)
    gt_captions, features, urls = minibatch
    gt_captions = decode_captions(gt_captions, data['idx_to_word'])

    sample_captions = small_lstm_model.sample(features)
    sample_captions = decode_captions(sample_captions, data['idx_to_word'])

    for gt_caption, sample_caption, url in zip(gt_captions, sample_captions, urls):
        plt.imshow(image_from_url(url))
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
        plt.axis('off')
        plt.show()
```

train
a man standing at home plate preparing to bat <END>
GT:<START> a man standing at home plate preparing to bat <END>



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



val
a dog <UNK> out on a bed under a blanket <END>
GT:<START> a man riding a skateboard <UNK> with <UNK> <UNK> on <END>

val
carrot cute cute dog standing on a motorcycle in a busy <END>
GT:<START> a person behind a stand with many oranges <END>

# Network Visualization (PyTorch)

In this notebook we will explore the use of *image gradients* for generating new images.

When training a model, we define a loss function which measures our current unhappiness with the model's performance; we then use backpropagation to compute the gradient of the loss with respect to the model parameters, and perform gradient descent on the model parameters to minimize the loss.

Here we will do something slightly different. We will start from a convolutional neural network model which has been pretrained to perform image classification on the ImageNet dataset. We will use this model to define a loss function which quantifies our current unhappiness with our image, then use backpropagation to compute the gradient of this loss with respect to the pixels of the image. We will then keep the model fixed, and perform gradient descent *on the image* to synthesize a new image which minimizes the loss.

In this notebook we will explore three techniques for image generation:

1. **Saliency Maps**: Saliency maps are a quick way to tell which part of the image influenced the classification decision made by the network.
2. **Fooling Images**: We can perturb an input image so that it appears the same to humans, but will be misclassified by the pretrained network.
3. **Class Visualization**: We can synthesize an image to maximize the classification score of a particular class; this can give us some sense of what the network is looking for when it classifies images of that class.

This notebook uses **PyTorch**; we have provided another notebook which explores the same concepts in TensorFlow. You only need to complete one of these two notebooks.

# Pretrained Model

For all of our image generation experiments, we will start with a convolutional neural network which was pretrained to perform image classification on ImageNet. We can use any model here, but for the purposes of this assignment we will use SqueezeNet [1], which achieves accuracies comparable to AlexNet but with a significantly reduced parameter count and computational complexity.

Using SqueezeNet rather than AlexNet or VGG or ResNet means that we can easily perform all image generation experiments on CPU.

[1] Iandola et al, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5MB model size", arXiv 2016

In [3]:

```python
# Download and load the pretrained SqueezeNet model.
model = torchvision.models.squeezenet1_1(pretrained=True)

# We don't want to train the model, so tell PyTorch not to compute gradients
# with respect to model parameters.
for param in model.parameters():
    param.requires_grad = False

# you may see warning regarding initialization deprecated, that's fine, please continue to next steps
```

# Saliency Maps

Using this pretrained model, we will compute class saliency maps as described in Section 3.1 of [2].

A **saliency map** tells us the degree to which each pixel in the image affects the classification score for that image. To compute it, we compute the gradient of the unnormalized score corresponding to the correct class (which is a scalar) with respect to the pixels of the image. If the image has shape `(3, H, W)` then this gradient will also have shape `(3, H, W)`; for each pixel in the image, this gradient tells us the amount by which the classification score will change if the pixel changes by a small amount. To compute the saliency map, we take the absolute value of this gradient, then take the maximum value over the 3 input channels; the final saliency map thus has shape `(H, W)` and all entries are nonnegative.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

```python
def compute_saliency_maps(X, y, model):
    """
    Compute a class saliency map using the model for images X and labels y.

    Input:
    - X: Input images; Tensor of shape (N, 3, H, W)
    - y: Labels for X; LongTensor of shape (N,)
    - model: A pretrained CNN that will be used to compute the saliency map.

    Returns:
    - saliency: A Tensor of shape (N, H, W) giving the saliency maps for the input
    images.
    """
    # Make sure the model is in "test" mode
    model.eval()

    # Make input tensor require gradient
    X.requires_grad_()

    saliency = None
    ##############################################################################
    # TODO: Implement this function. Perform a forward and backward pass through #
    # the model to compute the gradient of the correct class score with respect  #
    # to each input image. You first want to compute the loss over the correct   #
    # scores (we'll combine losses across a batch by summing), and then compute  #
    # the gradients with a backward pass.                                        #
    ##############################################################################
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    # Forward pass.
    scores = model(X)
    scores = scores.gather(1, y.view(-1, 1)).squeeze()

    # Backward pass, supply initial gradients of same tensor shape as scores.
    scores.backward(torch.ones(scores.size()))

    # Get gradient for image.
    saliency = X.grad

    # From 3d to 1d.
    saliency = saliency.abs()
    saliency, i = torch.max(saliency,dim=1)


    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    ##############################################################################
    #                              END OF YOUR CODE                              #
    ##############################################################################
    return saliency
```

Once you have completed the implementation in the cell above, run the following to visualize some class saliency maps on our example images from the ImageNet validation set:

```python
def show_saliency_maps(X, y):
    # Convert X and y from numpy arrays to Torch Tensors
    X_tensor = torch.cat([preprocess(Image.fromarray(x)) for x in X], dim=0)
    y_tensor = torch.LongTensor(y)

    # Compute saliency maps for images in X
    saliency = compute_saliency_maps(X_tensor, y_tensor, model)

    # Convert the saliency map from Torch Tensor to numpy array and show images
    # and saliency maps together.
    saliency = saliency.numpy()
    N = X.shape[0]
    for i in range(N):
        plt.subplot(2, N, i + 1)
        plt.imshow(X[i])
        plt.axis('off')
        plt.title(class_names[y[i]])
        plt.subplot(2, N, N + i + 1)
        plt.imshow(saliency[i], cmap=plt.cm.hot)
```
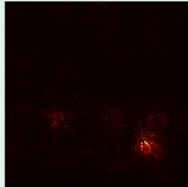
```
        plt.axis('off')
        plt.gcf().set_size_inches(12, 5)
    plt.show()

show_saliency_maps(X, y)
```
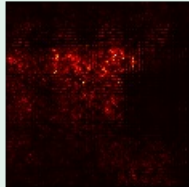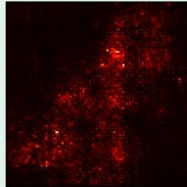
## INLINE QUESTION

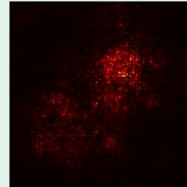A friend of yours suggests that in order to find an image that maximizes the correct score, we can perform gradient ascent on the input image, but instead of the gradient we can actually use the saliency map in each step to update the image. Is this assertion true? Why or why not?

**Your Answer:** No. For all 3 channels in input image, only one channel will be used for ascent. Some information may be lost.

# Fooling Images

We can also use image gradients to generate "fooling images" as discussed in [3]. Given an image and a target class, we can perform gradient **ascent** over the image to maximize the target class, stopping when the network classifies the image as the target class. Implement the following function to generate fooling images.

[3] Szegedy et al, "Intriguing properties of neural networks", ICLR 2014

In [7]:

```python
def make_fooling_image(X, target_y, model):
    """
    Generate a fooling image that is close to X, but that the model classifies
    as target_y.

    Inputs:
    - X: Input image; Tensor of shape (1, 3, 224, 224)
    - target_y: An integer in the range [0, 1000)
    - model: A pretrained CNN

    Returns:
    - X_fooling: An image that is close to X, but that is classifed as target_y
    by the model.
    """
    # Initialize our fooling image to the input image, and make it require gradient
    X_fooling = X.clone()
    X_fooling = X_fooling.requires_grad_()

    learning_rate = 1
    ##############################################################################
    # TODO: Generate a fooling image X_fooling that the model will classify as   #
    # the class target_y. You should perform gradient ascent on the score of the #
    # target class, stopping when the model is fooled.                           #
    # When computing an update step, first normalize the gradient:               #
    #   dX = learning_rate * g / ||g||_2                                         #
    #                                                                            #
    # You should write a training loop.                                          #
    #                                                                            #
```

```
    # HINT: For most examples, you should be able to generate a fooling image    #
    # in fewer than 100 iterations of gradient ascent.                           #
    # You can print your progress over iterations to check your algorithm.       #
    ##############################################################################
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    for i in range(100):
        scores = model(X_fooling)
        pred_idx = scores.data.max(dim=1)[1][0]
        if pred_idx == target_y:
            break
        target_score = scores[0, target_y]
        target_score.backward()

        # Gradient for image.
        grad = X_fooling.grad.data
        # Update the image with normalized gradient.
        X_fooling.data += learning_rate * (grad / grad.norm())
        X_fooling.grad.zero_()

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    ##############################################################################
    #                            END OF YOUR CODE                                #
    ##############################################################################
    return X_fooling
```

In [8]:

```
idx = 0
target_y = 6

X_tensor = torch.cat([preprocess(Image.fromarray(x)) for x in X], dim=0)
X_fooling = make_fooling_image(X_tensor[idx:idx+1], target_y, model)

scores = model(X_fooling)
assert target_y == scores.data.max(1)[1][0].item(), 'The model is not fooled!'
```

After generating a fooling image, run the following cell to visualize the original image, the fooling image, as well as the difference between them.

In [9]:

```
X_fooling_np = deprocess(X_fooling.clone())
X_fooling_np = np.asarray(X_fooling_np).astype(np.uint8)

plt.subplot(1, 4, 1)
plt.imshow(X[idx])
plt.title(class_names[y[idx]])
plt.axis('off')

plt.subplot(1, 4, 2)
plt.imshow(X_fooling_np)
plt.title(class_names[target_y])
plt.axis('off')

plt.subplot(1, 4, 3)
X_pre = preprocess(Image.fromarray(X[idx]))
diff = np.asarray(deprocess(X_fooling - X_pre, should_rescale=False))
plt.imshow(diff)
plt.title('Difference')
plt.axis('off')

plt.subplot(1, 4, 4)
diff = np.asarray(deprocess(10 * (X_fooling - X_pre), should_rescale=False))
plt.imshow(diff)
plt.title('Magnified difference (10x)')
plt.axis('off')

plt.gcf().set_size_inches(12, 5)
plt.show()
```
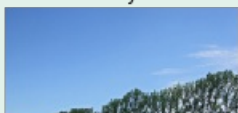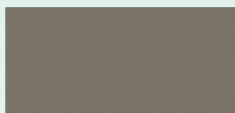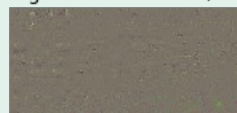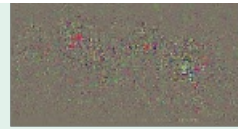
# Class visualization
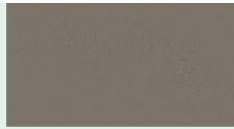
By starting with a random noise image and performing gradient ascent on a target class, we can generate an image that the network will recognize as the target class. This idea was first presented in [2]; [3] extended this idea by suggesting several regularization techniques that can improve the quality of the generated image.

Concretely, let $I$ be an image and let $y$ be a target class. Let $s_y(I)$ be the score that a convolutional network assigns to the image $I$ for class $y$; note that these are raw unnormalized scores, not class probabilities. We wish to generate an image $I^*$ that achieves a high score for the class $y$ by solving the problem

$$I^* = \arg\max_I (s_y(I) - R(I))$$

where $R$ is a (possibly implicit) regularizer (note the sign of $R(I)$ in the argmax: we want to minimize this regularization term). We can solve this optimization problem using gradient ascent, computing gradients with respect to the generated image. We will use (explicit) L2 regularization of the form

$$R(I) = \lambda \|I\|_2^2$$

**and** implicit regularization as suggested by [3] by periodically blurring the generated image. We can solve this problem using gradient ascent on the generated image.

In the cell below, complete the implementation of the `create_class_visualization` function.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

[3] Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML 2015 Deep Learning Workshop

In [13]:

```python
def create_class_visualization(target_y, model, dtype, **kwargs):
    """
    Generate an image to maximize the score of target_y under a pretrained model.

    Inputs:
    - target_y: Integer in the range [0, 1000) giving the index of the class
    - model: A pretrained CNN that will be used to generate the image
    - dtype: Torch datatype to use for computations

    Keyword arguments:
    - l2_reg: Strength of L2 regularization on the image
    - learning_rate: How big of a step to take
    - num_iterations: How many iterations to use
    - blur_every: How often to blur the image as an implicit regularizer
    - max_jitter: How much to gjitter the image as an implicit regularizer
    - show_every: How often to show the intermediate result
    """
    model.type(dtype)
    l2_reg = kwargs.pop('l2_reg', 1e-3)
    learning_rate = kwargs.pop('learning_rate', 25)
    num_iterations = kwargs.pop('num_iterations', 100)
    blur_every = kwargs.pop('blur_every', 10)
    max_jitter = kwargs.pop('max_jitter', 16)
    show_every = kwargs.pop('show_every', 25)

    # Randomly initialize the image as a PyTorch Tensor, and make it requires gradient.
    img = torch.randn(1, 3, 224, 224).mul_(1.0).type(dtype).requires_grad_()

    for t in range(num_iterations):
        # Randomly jitter the image a bit; this gives slightly nicer results
        ox, oy = random.randint(0, max_jitter), random.randint(0, max_jitter)
        img.data.copy_(jitter(img.data, ox, oy))

        ########################################################################
        # TODO: Use the model to compute the gradient of the score for the     #
```

```python
        # class target_y with respect to the pixels of the image, and make a   #
        # gradient step on the image using the learning rate. Don't forget the  #
        # L2 regularization term!                                               #
        # Be very careful about the signs of elements in your code.             #
        #########################################################################
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        scores = model(img)

        # Get the score for the target class.
        target_score = scores[0,target_y]

        # Backward pass to get gradient wrt image.
        target_score.backward()

        grad = img.grad.data - 2 * l2_reg * img

        img.data += learning_rate*grad
        img.grad.data.zero_()

        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
        #########################################################################
        #                          END OF YOUR CODE                             #
        #########################################################################

        # Undo the random jitter
        img.data.copy_(jitter(img.data, -ox, -oy))

        # As regularizer, clamp and periodically blur the image
        for c in range(3):
            lo = float(-SQUEEZENET_MEAN[c] / SQUEEZENET_STD[c])
            hi = float((1.0 - SQUEEZENET_MEAN[c]) / SQUEEZENET_STD[c])
            img.data[:, c].clamp_(min=lo, max=hi)
        if t % blur_every == 0:
            blur_image(img.data, sigma=0.5)

        # Periodically show the image
        if t == 0 or (t + 1) % show_every == 0 or t == num_iterations - 1:
            plt.imshow(deprocess(img.data.clone().cpu()))
            class_name = class_names[target_y]
            plt.title('%s\nIteration %d / %d' % (class_name, t + 1, num_iterations))
            plt.gcf().set_size_inches(4, 4)
            plt.axis('off')
            plt.show()

    return deprocess(img.data.cpu())
```

Once you have completed the implementation in the cell above, run the following cell to generate an image of a Tarantula:

In [14]:

```python
dtype = torch.FloatTensor
# dtype = torch.cuda.FloatTensor # Uncomment this to use GPU
model.type(dtype)

target_y = 76 # Tarantula
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
out = create_class_visualization(target_y, model, dtype)
```



```
# Backward pass to get gradient
```

tarantula
Iteration 25 / 100



tarantula
Iteration 50 / 100



tarantula
Iteration 75 / 100



tarantula
Iteration 100 / 100



Try out your class visualization on other classes! You should also feel free to play with various hyperparameters to try and improve the quality of the generated image, but this is not required.

In [15]:

```
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
target_y = np.random.randint(1000)
print(class_names[target_y])
X = create_class_visualization(target_y, model, dtype)
```

mountain tent

mountain tent
Iteration 1 / 100



mountain tent
Iteration 25 / 100



mountain tent
Iteration 50 / 100



mountain tent
Iteration 75 / 100

mountain tent
Iteration 100 / 100

In [ ]:

# Style Transfer

In this notebook we will implement the style transfer technique from ["Image Style Transfer Using Convolutional Neural Networks"](#) [(Gatys et al., CVPR 2015)](#).

The general idea is to take two images, and produce a new image that reflects the content of one but the artistic "style" of the other. We will do this by first formulating a loss function that matches the content and style of each respective image in the feature space of a deep network, and then performing gradient descent on the pixels of the image itself.

The deep network we use as a feature extractor is [SqueezeNet](#), a small model that has been trained on ImageNet. You could use any network, but we chose SqueezeNet here for its small size and efficiency.

Here's an example of the images you'll be able to produce by the end of this notebook:

## Computing Loss

We're going to compute the three components of our loss function now. The loss function is a weighted sum of three terms: content loss + style loss + total variation loss. You'll fill in the functions that compute these weighted terms below.

In [105]:

```python
def content_loss(content_weight, content_current, content_original):
    """
    Compute the content loss for style transfer.

    Inputs:
    - content_weight: Scalar giving the weighting for the content loss.
    - content_current: features of the current image; this is a PyTorch Tensor of shape
      (1, C_l, H_l, W_l).
    - content_target: features of the content image, Tensor with shape (1, C_l, H_l, W_l).

    Returns:
    - scalar content loss
    """
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

#     C_l, H_l, W_l = content_current.shape[1:]
#     F = content_current.reshape(C_l, H_l * W_l)
#     P = content_original.reshape(C_l, H_l * W_l)
#     Lc = content_weight * np.sum(np.square(F - P))

    Lc = content_weight * torch.sum((content_current - content_original) ** 2)
    return Lc

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

Test your content loss. You should see errors less than 0.001.

In [106]:

```python
def content_loss_test(correct):
    content_image = 'styles/tubingen.jpg'
    image_size =  192
    content_layer = 3
    content_weight = 6e-2

    c_feats, content_img_var = features_from_img(content_image, image_size)

    bad_img = torch.zeros(*content_img_var.data.size()).type(dtype)
    feats = extract_features(bad_img, cnn)

    student_output = content_loss(content_weight, c_feats[content_layer], feats[content_layer]).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Maximum error is {:.3f}'.format(error))
```

```
content_loss_test(answers['cl_out'])
```

Maximum error is 0.000

## Style loss

Now we can tackle the style loss. For a given layer $\ell$, the style loss is defined as follows:

First, compute the Gram matrix G which represents the correlations between the responses of each filter, where F is as above. The Gram matrix is an approximation to the covariance matrix -- we want the activation statistics of our generated image to match the activation statistics of our style image, and matching the (approximate) covariance is one way to do that. There are a variety of ways you could do this, but the Gram matrix is nice because it's easy to compute and in practice shows good results.

Given a feature map $F^\ell$ of shape $(C_\ell, M_\ell)$, the Gram matrix has shape $(C_\ell, C_\ell)$ and its elements are given by:

$$G_{ij}^\ell = \sum_k F_{ik}^\ell F_{jk}^\ell$$

Assuming $G^\ell$ is the Gram matrix from the feature map of the current image, $A^\ell$ is the Gram Matrix from the feature map of the source style image, and $w_\ell$ a scalar weight term, then the style loss for the layer $\ell$ is simply the weighted Euclidean distance between the two Gram matrices:

$$L_s^\ell = w_\ell \sum_{i,j} \left( G_{ij}^\ell - A_{ij}^\ell \right)^2$$

In practice we usually compute the style loss at a set of layers $L$ rather than just a single layer $\ell$; then the total style loss is the sum of style losses at each layer:

$$L_s = \sum_{\ell \in L} L_s^\ell$$

Begin by implementing the Gram matrix computation below:

In [107]:

```python
def gram_matrix(features, normalize=True):
    """
    Compute the Gram matrix from features.

    Inputs:
    - features: PyTorch Tensor of shape (N, C, H, W) giving features for
      a batch of N images.
    - normalize: optional, whether to normalize the Gram matrix
        If True, divide the Gram matrix by the number of neurons (H * W * C)

    Returns:
    - gram: PyTorch Tensor of shape (N, C, C) giving the
      (optionally normalized) Gram matrices for the N input images.
    """
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    N, C, H, W = features.shape
    features = features.reshape(N, C, H * W)
    gram = torch.zeros([N,C,C])

    for i in range(N):
        gram[i,:] = torch.mm(features[i,:], features[i,:].transpose(1, 0))

    if normalize == True:
        gram /= H * W * C * 1.0

    return gram


    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

Test your Gram matrix code. You should see errors less than 0.001.

In [108]:

```
def gram_matrix_test(correct):
    style_image = 'styles/starry_night.jpg'
    style_size = 192
    feats, _ = features_from_img(style_image, style_size)
    student_output = gram_matrix(feats[5].clone()).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Maximum error is {:.3f}'.format(error))

gram_matrix_test(answers['gm_out'])
```

```
Maximum error is 0.001
```

Next, implement the style loss:

In [109]:

```
# Now put it together in the style_loss function...
def style_loss(feats, style_layers, style_targets, style_weights):
    """
    Computes the style loss at a set of layers.

    Inputs:
    - feats: list of the features at every layer of the current image, as produced by
      the extract_features function.
    - style_layers: List of layer indices into feats giving the layers to include in the
      style loss.
    - style_targets: List of the same length as style_layers, where style_targets[i] is
      a PyTorch Tensor giving the Gram matrix of the source style image computed at
      layer style_layers[i].
    - style_weights: List of the same length as style_layers, where style_weights[i]
      is a scalar giving the weight for the style loss at layer style_layers[i].

    Returns:
    - style_loss: A PyTorch Tensor holding a scalar giving the style loss.
    """
    # Hint: you can do this with one for loop over the style layers, and should
    # not be very much code (~5 lines). You will need to use your gram_matrix function.
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    Ls = 0
    i = 0

    # Compute style loss for each desired feature layer and sum.
    for layer in style_layers:
        current = gram_matrix(feats[layer])
        Ls += style_weights[i] * torch.sum(torch.pow((current - style_targets[i]), 2))
        i += 1

    return Ls

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

Test your style loss implementation. The error should be less than 0.001.

In [110]:

```
def style_loss_test(correct):
    content_image = 'styles/tubingen.jpg'
    style_image = 'styles/starry_night.jpg'
    image_size =  192
    style_size = 192
    style_layers = [1, 4, 6, 7]
    style_weights = [300000, 1000, 15, 3]

    c_feats, _ = features_from_img(content_image, image_size)
    feats, _ = features_from_img(style_image, style_size)
    style_targets = []
    for idx in style_layers:
        style_targets.append(gram_matrix(feats[idx].clone()))

    student_output = style_loss(c_feats, style_layers, style_targets, style_weights).cpu().data.num
py()
    error = rel_error(correct, student_output)
    print('Error is {:.3f}'.format(error))
```

```
style_loss_test(answers['sl_out'])
```

Error is 0.000

## Total-variation regularization

It turns out that it's helpful to also encourage smoothness in the image. We can do this by adding another term to our loss that penalizes wiggles or "total variation" in the pixel values.

You can compute the "total variation" as the sum of the squares of differences in the pixel values for all pairs of pixels that are next to each other (horizontally or vertically). Here we sum the total-variation regualarization for each of the 3 input channels (RGB), and weight the total summed loss by the total variation weight, $w_t$:

$$L_{tv} = w_t \times \left( \sum_{c=1}^{3}\sum_{i=1}^{H-1}\sum_{j=1}^{W}(x_{i+1,j,c} - x_{i,j,c})^2 + \sum_{c=1}^{3}\sum_{i=1}^{H}\sum_{j=1}^{W-1}(x_{i,j+1,c} - x_{i,j,c})^2 \right)$$

In the next cell, fill in the definition for the TV loss term. To receive full credit, your implementation should not have any loops.

In [111]:

```
def tv_loss(img, tv_weight):
    """
    Compute total variation loss.

    Inputs:
    - img: PyTorch Variable of shape (1, 3, H, W) holding an input image.
    - tv_weight: Scalar giving the weight w_t to use for the TV loss.

    Returns:
    - loss: PyTorch Variable holding a scalar giving the total variation loss
      for img weighted by tv_weight.
    """
    # Your implementation should be vectorized and not require any loops!
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    Ltv = 0
    H_direction = torch.sum((img[:, :, 1:, :] - img[:, :, :-1, :])**2)
    W_direction = torch.sum((img[:, :, :, 1:] - img[:, :, :, :-1])**2)
    Ltv = tv_weight * (H_direction + W_direction)
    return Ltv

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

Test your TV loss implementation. Error should be less than 0.0001.

In [112]:

```
def tv_loss_test(correct):
    content_image = 'styles/tubingen.jpg'
    image_size =  192
    tv_weight = 2e-2

    content_img = preprocess(PIL.Image.open(content_image), size=image_size).type(dtype)

    student_output = tv_loss(content_img, tv_weight).cpu().data.numpy()
    error = rel_error(correct, student_output)
    print('Error is {:.3f}'.format(error))

tv_loss_test(answers['tv_out'])
```

Error is 0.000

Now we're ready to string it all together (you shouldn't have to modify this function):

## Generate some pretty pictures!

Try out `style_transfer` on the three different parameter sets below. Make sure to run all three cells. Feel free to add your own, but make sure to include the results of style transfer on the third parameter set (starry night) in your submitted notebook.

- The `content_image` is the filename of content image.
- The `style_image` is the filename of style image.

- The `image_size` is the size of smallest image dimension of the content image (used for content loss and generated image).
- The `style_size` is the size of smallest style image dimension.
- The `content_layer` specifies which layer to use for content loss.
- The `content_weight` gives weighting on content loss in the overall loss function. Increasing the value of this parameter will make the final image look more realistic (closer to the original content).
- `style_layers` specifies a list of which layers to use for style loss.
- `style_weights` specifies a list of weights to use for each layer in style_layers (each of which will contribute a term to the overall style loss). We generally use higher weights for the earlier style layers because they describe more local/smaller scale features, which are more important to texture than features over larger receptive fields. In general, increasing these weights will make the resulting image look less like the original content and more distorted towards the appearance of the style image.
- `tv_weight` specifies the weighting of total variation regularization in the overall loss function. Increasing this value makes the resulting image look smoother and less jagged, at the cost of lower fidelity to style and content.

Below the next three cells of code (in which you shouldn't change the hyperparameters), feel free to copy and paste the parameters to play around them and see how the resulting image changes.

In [114]:

```python
# Composition VII + Tubingen
params1 = {
    'content_image' : 'styles/tubingen.jpg',
    'style_image' : 'styles/composition_vii.jpg',
    'image_size' : 192,
    'style_size' : 512,
    'content_layer' : 3,
    'content_weight' : 5e-2,
    'style_layers' : (1, 4, 6, 7),
    'style_weights' : (20000, 500, 12, 1),
    'tv_weight' : 5e-2
}

style_transfer(**params1)
```

Content Source Img.    Style Source Img.



Iteration 0



Iteration 100

Iteration 199

```python
# Scream + Tubingen
params2 = {
    'content_image':'styles/tubingen.jpg',
    'style_image':'styles/the_scream.jpg',
    'image_size':192,
    'style_size':224,
    'content_layer':3,
    'content_weight':3e-2,
    'style_layers':[1, 4, 6, 7],
    'style_weights':[200000, 800, 12, 1],
    'tv_weight':2e-2
}

style_transfer(**params2)
```

Content Source Img.

Style Source Img.



Iteration 0



Iteration 100

Iteration 199

```python
# Starry Night + Tubingen
params3 = {
    'content_image' : 'styles/tubingen.jpg',
    'style_image' : 'styles/starry_night.jpg',
    'image_size' : 192,
    'style_size' : 192,
    'content_layer' : 3,
    'content_weight' : 6e-2,
    'style_layers' : [1, 4, 6, 7],
    'style_weights' : [300000, 1000, 15, 3],
    'tv_weight' : 2e-2
}

style_transfer(**params3)
```



Iteration 0



Iteration 100

```
Iteration 199
```



# Feature Inversion

The code you've written can do another cool thing. In an attempt to understand the types of features that convolutional networks learn to recognize, a recent paper [1] attempts to reconstruct an image from its feature representation. We can easily implement this idea using image gradients from the pretrained network, which is exactly what we did above (but with two different feature representations).

Now, if you set the style weights to all be 0 and initialize the starting image to random noise instead of the content source image, you'll reconstruct an image from the feature representation of the content source image. You're starting with total noise, but you should end up with something that looks quite a bit like your original image.

(Similarly, you could do "texture synthesis" from scratch if you set the content weight to 0 and initialize the starting image to random noise, but we won't ask you to do that here.)

Run the following cell to try out feature inversion.

[1] Aravindh Mahendran, Andrea Vedaldi, "Understanding Deep Image Representations by Inverting them", CVPR 2015

In [117]:

```python
# Feature Inversion -- Starry Night + Tubingen
params_inv = {
    'content_image' : 'styles/tubingen.jpg',
    'style_image' : 'styles/starry_night.jpg',
    'image_size' : 192,
    'style_size' : 192,
    'content_layer' : 3,
    'content_weight' : 6e-2,
    'style_layers' : [1, 4, 6, 7],
    'style_weights' : [0, 0, 0, 0], # we discard any contributions from style to the loss
    'tv_weight' : 2e-2,
    'init_random': True # we want to initialize our image to be random
}

style_transfer(**params_inv)
```



Content Source Img.        Style Source Img.

Iteration 0



Iteration 100



Iteration 199

# Generative Adversarial Networks (GANs)

So far in CS231N, all the applications of neural networks that we have explored have been **discriminative models** that take an input and are trained to produce a labeled output. This has ranged from straightforward classification of image categories to sentence generation (which was still phrased as a classification problem, our labels were in vocabulary space and we'd learned a recurrence to capture multi-word labels). In this notebook, we will expand our repetoire, and build **generative models** using neural networks. Specifically, we will learn how to build models which generate novel images that resemble a set of training images.

## What is a GAN?

In 2014, [Goodfellow et al.](#) presented a method for training generative models called Generative Adversarial Networks (GANs for short). In a GAN, we build two different neural networks. Our first network is a traditional classification network, called the **discriminator**. We will train the discriminator to take images, and classify them as being real (belonging to the training set) or fake (not present in the training set). Our other network, called the **generator**, will take random noise as input and transform it using a neural network to produce images. The goal of the generator is to fool the discriminator into thinking the images it produced are real.

We can think of this back and forth process of the generator ($G$) trying to fool the discriminator ($D$), and the discriminator trying to correctly classify real vs. fake as a minimax game:

$$\underset{G}{\text{minimize}}\, \underset{D}{\text{maximize}}\, E_{x \sim p_{\text{data}}}[\log D(x)] + E_{z \sim p(z)}[\log(1 - D(G(z)))]$$

where $z \sim p(z)$ are the random noise samples, $G(z)$ are the generated images using the neural network generator $G$, and $D$ is the output of the discriminator, specifying the probability of an input being real. In [Goodfellow et al.](#), they analyze this minimax game and show how it relates to minimizing the Jensen-Shannon divergence between the training data distribution and the generated samples from $G$.

To optimize this minimax game, we will alternate between taking gradient *descent* steps on the objective for $G$, and gradient *ascent* steps on the objective for $D$:

1. update the **generator** ($G$) to minimize the probability of the **discriminator making the correct choice**.
2. update the **discriminator** ($D$) to maximize the probability of the **discriminator making the correct choice**.

While these updates are useful for analysis, they do not perform well in practice. Instead, we will use a different objective when we update the generator: maximize the probability of the **discriminator making the incorrect choice**. This small change helps to alleviaite problems with the generator gradient vanishing when the discriminator is confident. This is the standard update used in most GAN papers, and was used in the original paper from [Goodfellow et al.](#).

In this assignment, we will alternate the following updates:

1. Update the generator ($G$) to maximize the probability of the discriminator making the incorrect choice on generated data:

$$\underset{G}{\text{maximize}}\, E_{z \sim p(z)}[\log D(G(z))]$$

2. Update the discriminator ($D$), to maximize the probability of the discriminator making the correct choice on real and generated data:

$$\underset{D}{\text{maximize}}\, E_{x \sim p_{\text{data}}}[\log D(x)] + E_{z \sim p(z)}[\log(1 - D(G(z)))]$$

## What else is there?

Since 2014, GANs have exploded into a huge research area, with massive [workshops](#), and [hundreds of new papers](#). Compared to other approaches for generative models, they often produce the highest quality samples but are some of the most difficult and finicky models to train (see [this github repo](#) that contains a set of 17 hacks that are useful for getting models working). Improving the stabiilty and robustness of GAN training is an open research question, with new papers coming out every day! For a more recent tutorial on GANs, see [here](#). There is also some even more recent exciting work that changes the objective function to Wasserstein distance and yields much more stable results across model architectures: [WGAN](#), [WGAN-GP](#).

GANs are not the only way to train a generative model! For other approaches to generative modeling check out the [deep generative model chapter](#) of the Deep Learning [book](#). Another popular way of training neural networks as generative models is Variational Autoencoders (co-discovered [here](#) and [here](#)). Variatonal autoencoders combine neural networks with variationl inference to train deep generative models. These models tend to be far more stable and easier to train but currently don't produce samples that are

as pretty as GANs.

Here's an example of what your outputs from the 3 different models you're going to train should look like... note that GANs are sometimes finicky, so your outputs might not look *exactly* like this... this is just meant to be a *rough* guideline of the kind of quality you can expect:

## Random Noise

Generate uniform noise from -1 to 1 with shape `[batch_size, dim]`.

Hint: use `torch.rand`.

In [4]:

```python
def sample_noise(batch_size, dim):
    """
    Generate a PyTorch Tensor of uniform random noise.

    Input:
    - batch_size: Integer giving the batch size of noise to generate.
    - dim: Integer giving the dimension of noise to generate.

    Output:
    - A PyTorch Tensor of shape (batch_size, dim) containing uniform
      random noise in the range (-1, 1).
    """
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    return torch.rand(batch_size, dim) * 2 - 1

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

Make sure noise is the correct shape and type:

In [5]:

```python
def test_sample_noise():
    batch_size = 3
    dim = 4
    torch.manual_seed(231)
    z = sample_noise(batch_size, dim)
    np_z = z.cpu().numpy()
    assert np_z.shape == (batch_size, dim)
    assert torch.is_tensor(z)
    assert np.all(np_z >= -1.0) and np.all(np_z <= 1.0)
    assert np.any(np_z < 0.0) and np.any(np_z > 0.0)
    print('All tests passed!')

test_sample_noise()
```

```
All tests passed!
```

## Discriminator

Our first step is to build a discriminator. Fill in the architecture as part of the `nn.Sequential` constructor in the function below. All fully connected layers should include bias terms. The architecture is:

- Fully connected layer with input size 784 and output size 256
- LeakyReLU with alpha 0.01
- Fully connected layer with input_size 256 and output size 256
- LeakyReLU with alpha 0.01
- Fully connected layer with input size 256 and output size 1

Recall that the Leaky ReLU nonlinearity computes $f(x) = \max$ for some fixed constant \alpha; for the LeakyReLU nonlinearities in the architecture above we set \alpha=0.01.

The output of the discriminator should have shape `[batch_size, 1]`, and contain real numbers corresponding to the scores that each of the `batch_size` inputs is a real image.

```python
def discriminator():
    """
    Build and return a PyTorch model implementing the architecture above.
    """
    model = nn.Sequential(
            # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
            Flatten(),
            nn.Linear(784, 256),
            nn.LeakyReLU(0.01),
            nn.Linear(256, 256),
            nn.LeakyReLU(0.01),
            nn.Linear(256, 1)

            # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    )
    return model
```

Test to make sure the number of parameters in the discriminator is correct:

In [9]:

```python
def test_discriminator(true_count=267009):
    model = discriminator()
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in discriminator. Check your achitecture.')
    else:
        print('Correct number of parameters in discriminator.')

test_discriminator()
```

```
Correct number of parameters in discriminator.
```

## Generator

Now to build the generator network:

- Fully connected layer from noise_dim to 1024
- `ReLU`
- Fully connected layer with size 1024
- `ReLU`
- Fully connected layer with size 784
- `TanH` (to clip the image to be in the range of [-1,1])

In [10]:

```python
def generator(noise_dim=NOISE_DIM):
    """
    Build and return a PyTorch model implementing the architecture above.
    """
    model = nn.Sequential(
            # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

            nn.Linear(noise_dim, 1024),
            nn.ReLU(),
            nn.Linear(1024, 1024),
            nn.ReLU(),
            nn.Linear(1024, 784),
            nn.Tanh()

            # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    )
    return model
```

Test to make sure the number of parameters in the generator is correct:

In [11]:

```python
def test_generator(true_count=1858320):
    model = generator(4)
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in generator. Check your achitecture.')
    else:
        print('Correct number of parameters in generator.')

test_generator()
```

```
Correct number of parameters in generator.
```

## GAN Loss

Compute the generator and discriminator loss. The generator loss is: $\ell_G = -\mathbb{E}_{z \sim p(z)}\left[\log D(G(z))\right]$ and the discriminator loss is: $\ell_D = -\mathbb{E}_{x \sim p_\text{data}}\left[\log D(x)\right] - \mathbb{E}_{z \sim p(z)}\left[\log \left(1-D(G(z))\right)\right]$ Note that these are negated from the equations presented earlier as we will be *minimizing* these losses.

**HINTS**: You should use the `bce_loss` function defined below to compute the binary cross entropy loss which is needed to compute the log probability of the true label given the logits output from the discriminator. Given a score $s\in\mathbb{R}$ and a label $y\in\{0, 1\}$, the binary cross entropy loss is
bce(s, y) = -y * \log(s) - (1 - y) * \log(1 - s)

A naive implementation of this formula can be numerically unstable, so we have provided a numerically stable implementation for you below.

You will also need to compute labels corresponding to real or fake and use the logit arguments to determine their size. Make sure you cast these labels to the correct data type using the global `dtype` variable, for example:

```
true_labels = torch.ones(size).type(dtype)
```

Instead of computing the expectation of $\log D(G(z))$, $\log D(x)$ and $\log \left(1-D(G(z))\right)$, we will be averaging over elements of the minibatch, so make sure to combine the loss by averaging instead of summing.

In [26]:

```python
def discriminator_loss(logits_real, logits_fake):
    """
    Computes the discriminator loss described above.

    Inputs:
    - logits_real: PyTorch Tensor of shape (N,) giving scores for the real data.
    - logits_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Returns:
    - loss: PyTorch Tensor containing (scalar) the loss for the discriminator.
    """
    N = logits_real.size()
    true_labels = torch.ones(N).type(dtype)
    loss_true = bce_loss(logits_real, true_labels)

    fake_labels = 1 - true_labels
    loss_fake = bce_loss(logits_fake, fake_labels)

    loss = loss_true + loss_fake
    return loss


def generator_loss(logits_fake):
    """
    Computes the generator loss described above.

    Inputs:
    - logits_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Returns:
    - loss: PyTorch Tensor containing the (scalar) loss for the generator.
    """
    N = logits_fake.size()
    true_labels = torch.ones(N).type(dtype)
    loss = bce_loss(logits_fake, true_labels)
    return loss
```

Test your generator and discriminator loss. You should see errors < 1e-7.

```python
def test_discriminator_loss(logits_real, logits_fake, d_loss_true):
    d_loss = discriminator_loss(torch.Tensor(logits_real).type(dtype),
                                torch.Tensor(logits_fake).type(dtype)).cpu().numpy()
    print("Maximum error in d_loss: %g"%rel_error(d_loss_true, d_loss))

test_discriminator_loss(answers['logits_real'], answers['logits_fake'],
                        answers['d_loss_true'])
```

```
Maximum error in d_loss: 2.83811e-08
```

```python
def test_generator_loss(logits_fake, g_loss_true):
    g_loss = generator_loss(torch.Tensor(logits_fake).type(dtype)).cpu().numpy()
    print("Maximum error in g_loss: %g"%rel_error(g_loss_true, g_loss))

test_generator_loss(answers['logits_fake'], answers['g_loss_true'])
```

```
Maximum error in g_loss: 4.4518e-09
```

# Optimizing our loss

Make a function that returns an `optim.Adam` optimizer for the given model with a 1e-3 learning rate, beta1=0.5, beta2=0.999. You'll use this to construct optimizers for the generators and discriminators for the rest of the notebook.

```python
def get_optimizer(model):
    """
    Construct and return an Adam optimizer for the model with learning rate 1e-3,
    beta1=0.5, and beta2=0.999.

    Input:
    - model: A PyTorch model that we want to optimize.

    Returns:
    - An Adam optimizer for the model with the desired hyperparameters.
    """
    optimizer = optim.Adam(model.parameters(), lr = 1e-3, betas = (0.5,0.999))
    return optimizer
```

# Training a GAN!

We provide you the main training loop... you won't need to change this function, but we encourage you to read through and understand it.

```python
# Make the discriminator
D = discriminator().type(dtype)

# Make the generator
G = generator().type(dtype)

# Use the function you wrote earlier to get optimizers for the Discriminator and the Generator
D_solver = get_optimizer(D)
G_solver = get_optimizer(G)
# Run it!
run_a_gan(D, G, D_solver, G_solver, discriminator_loss, generator_loss)
```

```
Iter: 0, D: 1.328, G:0.7202
```

Iter: 250, D: 1.939, G:0.8623



Iter: 500, D: 1.212, G:1.01



Iter: 750, D: 1.407, G:0.7779



Iter: 1000, D: 1.138, G:1.107



Iter: 1250, D: 1.201, G:0.9872

Iter: 1500, D: 1.204, G:0.9861



Iter: 1750, D: 1.228, G:0.911



Iter: 2000, D: 1.188, G:0.7323



Iter: 2250, D: 1.267, G:0.9026



Iter: 2500, D: 1.444, G:0.7764

Iter: 2750, D: 1.282, G:0.8262



Iter: 3000, D: 1.287, G:0.9628



Iter: 3250, D: 1.312, G:0.8871



Iter: 3500, D: 1.273, G:0.8278



Iter: 3750, D: 1.273, G:0.8697

## Least Squares GAN

We'll now look at [Least Squares GAN](#), a newer, more stable alernative to the original GAN loss function. For this part, all we have to do is change the loss function and retrain the model. We'll implement equation (9) in the paper, with the generator loss: $\ell_G = \frac{1}{2}\mathbb{E}_{z \sim p(z)}\left[\left(D(G(z))-1\right)^2\right]$ and the discriminator loss: $\ell_D = \frac{1}{2}\mathbb{E}_{x \sim p_\text{data}}\left[\left(D(x)-1\right)^2\right] + \frac{1}{2}\mathbb{E}_{z \sim p(z)}\left[ \left(D(G(z))\right)^2\right]$

**HINTS**: Instead of computing the expectation, we will be averaging over elements of the minibatch, so make sure to combine the loss by averaging instead of summing. When plugging in for $D(x)$ and $D(G(z))$ use the direct output from the discriminator ( `scores_real` and `scores_fake` ).

In [17]:

```python
def ls_discriminator_loss(scores_real, scores_fake):
    """
    Compute the Least-Squares GAN loss for the discriminator.

    Inputs:
    - scores_real: PyTorch Tensor of shape (N,) giving scores for the real data.
    - scores_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Outputs:
    - loss: A PyTorch Tensor containing the loss.
    """
    loss = None
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    N = scores_real.size()
    loss = (0.5 * torch.mean((scores_real - 1)**2)) + (0.5 * torch.mean(scores_fake**2))

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    return loss

def ls_generator_loss(scores_fake):
    """
    Computes the Least-Squares GAN loss for the generator.

    Inputs:
    - scores_fake: PyTorch Tensor of shape (N,) giving scores for the fake data.

    Outputs:
    - loss: A PyTorch Tensor containing the loss.
    """
    loss = 0.5 * torch.mean((scores_fake - 1) ** 2)
    return loss
```

Before running a GAN with our new loss function, let's check it:

In [18]:

```python
def test_lsgan_loss(score_real, score_fake, d_loss_true, g_loss_true):
    score_real = torch.Tensor(score_real).type(dtype)
    score_fake = torch.Tensor(score_fake).type(dtype)
    d_loss = ls_discriminator_loss(score_real, score_fake).cpu().numpy()
    g_loss = ls_generator_loss(score_fake).cpu().numpy()
    print("Maximum error in d_loss: %g"%rel_error(d_loss_true, d_loss))
    print("Maximum error in g_loss: %g"%rel_error(g_loss_true, g_loss))
```

```
test_lsgan_loss(answers['logits_real'], answers['logits_fake'],
                answers['d_loss_lsgan_true'], answers['g_loss_lsgan_true'])
```

```
Maximum error in d_loss: 1.64377e-08
Maximum error in g_loss: 3.36961e-08
```

Run the following cell to train your model!

In [22]:

```
D_LS = discriminator().type(dtype)
G_LS = generator().type(dtype)

D_LS_solver = get_optimizer(D_LS)
G_LS_solver = get_optimizer(G_LS)

run_a_gan(D_LS, G_LS, D_LS_solver, G_LS_solver, ls_discriminator_loss, ls_generator_loss)
```

```
Iter: 0, D: 0.476, G:0.4826
```



```
Iter: 250, D: 0.1693, G:0.3271
```



```
Iter: 500, D: 0.1771, G:0.8023
```



```
Iter: 750, D: 0.1909, G:0.1692
```

Iter: 1000, D: 0.1256, G:0.3286



Iter: 1250, D: 0.1562, G:0.2851



Iter: 1500, D: 0.2746, G:0.12



Iter: 1750, D: 0.221, G:0.1799



Iter: 2000, D: 0.2323, G:0.1965

Iter: 2250, D: 0.2279, G:0.1575



Iter: 2500, D: 0.2109, G:0.1929



Iter: 2750, D: 0.2293, G:0.158



Iter: 3000, D: 0.2388, G:0.1622



Iter: 3250, D: 0.2377, G:0.1619

```
Iter: 3500, D: 0.2396, G:0.1667
```



```
Iter: 3750, D: 0.218, G:0.1694
```



# Deeply Convolutional GANs

In the first part of the notebook, we implemented an almost direct copy of the original GAN network from Ian Goodfellow. However, this network architecture allows no real spatial reasoning. It is unable to reason about things like "sharp edges" in general because it lacks any convolutional layers. Thus, in this section, we will implement some of the ideas from DCGAN, where we use convolutional networks

**Discriminator**

We will use a discriminator inspired by the TensorFlow MNIST classification tutorial, which is able to get above 99% accuracy on the MNIST dataset fairly quickly.

- Reshape into image tensor (Use Unflatten!)
- Conv2D: 32 Filters, 5x5, Stride 1
- Leaky ReLU(alpha=0.01)
- Max Pool 2x2, Stride 2
- Conv2D: 64 Filters, 5x5, Stride 1
- Leaky ReLU(alpha=0.01)
- Max Pool 2x2, Stride 2
- Flatten
- Fully Connected with output size 4 x 4 x 64
- Leaky ReLU(alpha=0.01)
- Fully Connected with output size 1

In [91]:

```python
def build_dc_classifier():
    """
    Build and return a PyTorch model for the DCGAN discriminator implementing
```

```
        the architecture above.
        """
    return nn.Sequential(
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        Unflatten(batch_size, 1, 28, 28),
        nn.Conv2d(1, 32, kernel_size = 5, stride = 1),
        nn.LeakyReLU(0.01),
        nn.MaxPool2d(kernel_size = 2, stride = 2),
        nn.Conv2d(32, 64, kernel_size = 5, stride = 1),
        nn.LeakyReLU(0.01),
        nn.MaxPool2d(kernel_size = 2, stride = 2),
        Flatten(),
        nn.Linear(4*4*64, 4*4*64),
        nn.LeakyReLU(0.01),
        nn.Linear(4*4*64, 1)

        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    )

data = next(enumerate(loader_train))[-1][0].type(dtype)
b = build_dc_classifier().type(dtype)
out = b(data)
print(out.size())
```

```
torch.Size([128, 1])
```

Check the number of parameters in your classifier as a sanity check:

In [92]:

```
def test_dc_classifer(true_count=1102721):
    model = build_dc_classifier()
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in generator. Check your achitecture.')
    else:
        print('Correct number of parameters in generator.')

test_dc_classifer()
```

```
Correct number of parameters in generator.
```

**Generator**

For the generator, we will copy the architecture exactly from the InfoGAN paper. See Appendix C.1 MNIST. See the documentation for tf.nn.conv2d_transpose. We are always "training" in GAN mode.

- Fully connected with output size 1024
- `ReLU`
- BatchNorm
- Fully connected with output size 7 x 7 x 128
- ReLU
- BatchNorm
- Reshape into Image Tensor of shape 7, 7, 128
- Conv2D^T (Transpose): 64 filters of 4x4, stride 2, 'same' padding (use `padding=1`)
- `ReLU`
- BatchNorm
- Conv2D^T (Transpose): 1 filter of 4x4, stride 2, 'same' padding (use `padding=1`)
- `TanH`
- Should have a 28x28x1 image, reshape back into 784 vector

In [93]:

```
def build_dc_generator(noise_dim=NOISE_DIM):
    """
    Build and return a PyTorch model implementing the DCGAN generator using
    the architecture described above.
    """
    return nn.Sequential(
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

```
        nn.Linear(noise_dim, 1024),
        nn.ReLU(),
        nn.BatchNorm1d(1024),
        nn.Linear(1024, 7*7*128),
        nn.ReLU(),
        nn.BatchNorm1d(7*7*128),
        Unflatten(batch_size, 128, 7, 7),
        nn.ConvTranspose2d(128, 64, kernel_size = 4, stride = 2, padding = 1),
        nn.ReLU(),
        nn.BatchNorm2d(64),
        nn.ConvTranspose2d(64, 1, kernel_size = 4, stride = 2, padding = 1),
        nn.Tanh(),
        Flatten()

        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    )

test_g_gan = build_dc_generator().type(dtype)
test_g_gan.apply(initialize_weights)

fake_seed = torch.randn(batch_size, NOISE_DIM).type(dtype)
fake_images = test_g_gan.forward(fake_seed)
fake_images.size()
```

Out[93]:

```
torch.Size([128, 784])
```

Check the number of parameters in your generator as a sanity check:

In [94]:

```
def test_dc_generator(true_count=6580801):
    model = build_dc_generator(4)
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in generator. Check your achitecture.')
    else:
        print('Correct number of parameters in generator.')

test_dc_generator()
```

```
Correct number of parameters in generator.
```

In [95]:

```
D_DC = build_dc_classifier().type(dtype)
D_DC.apply(initialize_weights)
G_DC = build_dc_generator().type(dtype)
G_DC.apply(initialize_weights)

D_DC_solver = get_optimizer(D_DC)
G_DC_solver = get_optimizer(G_DC)

run_a_gan(D_DC, G_DC, D_DC_solver, G_DC_solver, discriminator_loss, generator_loss, num_epochs=5)
```

```
Iter: 0, D: 1.541, G:0.5609
```



```
Iter: 250, D: 1.448, G:0.2154
```

Iter: 500, D: 1.349, G:1.197



Iter: 750, D: 1.273, G:1.437



Iter: 1000, D: 1.214, G:1.03



Iter: 1250, D: 1.208, G:1.206



Iter: 1500, D: 1.102, G:0.9956

```
Iter: 1750, D: 1.223, G:1.101
```



## INLINE QUESTION 1

We will look at an example to see why alternating minimization of the same objective (like in a GAN) can be tricky business.

Consider $f(x,y)=xy$. What does $\min_x\max_y f(x,y)$ evaluate to? (Hint: minmax tries to minimize the maximum value achievable.)

Now try to evaluate this function numerically for 6 steps, starting at the point $(1,1)$, by using alternating gradient (first updating y, then updating x using that updated y) with step size $1$. **Here step size is the learning_rate, and steps will be learning_rate * gradient.** You'll find that writing out the update step in terms of $x\_t, y\_t, x\_{t+1}, y\_{t+1}$ will be useful.

Breifly explain what $\min_x\max_y f(x,y)$ evaluates to and record the six pairs of explicit values for $(x\_t, y\_t)$ in the table below.

**Your answer:**

| y_0 | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-----|-----|-----|-----|-----|-----|-----|
| 1   | 2   | 0   | 0   | 0   | 0   | 0   |

| x_0 | x_1 | x_2 | x_3 | x_4 | x_5 | x_6 |
|-----|-----|-----|-----|-----|-----|-----|
| 1   | -2  | 0   | 0   | 0   | 0   | 0   |

$\min_x\max_y f(x, y)$ evaluates to 0. As we can see from the table above, within 6 steps we get the answer.

## INLINE QUESTION 2

Using this method, will we ever reach the optimal value? Why or why not?

**Your answer:**

Yes. We can see we got the same answers after y2, after two steps we reach the optimal value and remain unchanged.

## INLINE QUESTION 3

If the generator loss decreases during training while the discriminator loss stays at a constant high value from the start, is this a good sign? Why or why not? A qualitative answer is sufficient.

**Your answer:**

It is not a good sign. If generator loss decreases and discriminator loss stays at a constant high value, it means discriminator can't catch up with generator. In this case, discriminator is decreasing the probability of real data being considered to be real while generator is generating data. But we actually want them to reach an equilibrium. So this is not a good sign.

In [ ]:

# 1 rnn.py

```python
from builtins import range
from builtins import object
import numpy as np

from cs231n.layers import *
from cs231n.rnn_layers import *


class CaptioningRNN(object):
    """
    A CaptioningRNN produces captions from image features using a recurrent
    neural network.

    The RNN receives input vectors of size D, has a vocab size of V, works on
    sequences of length T, has an RNN hidden dimension of H, uses word vectors
    of dimension W, and operates on minibatches of size N.

    Note that we don't use any regularization for the CaptioningRNN.
    """

    def __init__(self, word_to_idx, input_dim=512, wordvec_dim=128,
                 hidden_dim=128, cell_type='rnn', dtype=np.float32):
        """
        Construct a new CaptioningRNN instance.

        Inputs:
        - word_to_idx: A dictionary giving the vocabulary. It contains V entries,
          and maps each string to a unique integer in the range [0, V).
        - input_dim: Dimension D of input image feature vectors.
        - wordvec_dim: Dimension W of word vectors.
        - hidden_dim: Dimension H for the hidden state of the RNN.
        - cell_type: What type of RNN to use; either 'rnn' or 'lstm'.
        - dtype: numpy datatype to use; use float32 for training and float64 for
          numeric gradient checking.
        """
        if cell_type not in {'rnn', 'lstm'}:
            raise ValueError('Invalid cell_type "%s"' % cell_type)

        self.cell_type = cell_type
        self.dtype = dtype
        self.word_to_idx = word_to_idx
        self.idx_to_word = {i: w for w, i in word_to_idx.items()}
        self.params = {}

        vocab_size = len(word_to_idx)

        self._null = word_to_idx['<NULL>']
        self._start = word_to_idx.get('<START>', None)
        self._end = word_to_idx.get('<END>', None)

        # Initialize word vectors
        self.params['W_embed'] = np.random.randn(vocab_size, wordvec_dim)
        self.params['W_embed'] /= 100

        # Initialize CNN -> hidden state projection parameters
        self.params['W_proj'] = np.random.randn(input_dim, hidden_dim)
        self.params['W_proj'] /= np.sqrt(input_dim)
        self.params['b_proj'] = np.zeros(hidden_dim)

        # Initialize parameters for the RNN
        dim_mul = {'lstm': 4, 'rnn': 1}[cell_type]
        self.params['Wx'] = np.random.randn(wordvec_dim, dim_mul * hidden_dim)
        self.params['Wx'] /= np.sqrt(wordvec_dim)
        self.params['Wh'] = np.random.randn(hidden_dim, dim_mul * hidden_dim)
        self.params['Wh'] /= np.sqrt(hidden_dim)
        self.params['b'] = np.zeros(dim_mul * hidden_dim)

        # Initialize output to vocab weights
        self.params['W_vocab'] = np.random.randn(hidden_dim, vocab_size)
        self.params['W_vocab'] /= np.sqrt(hidden_dim)
        self.params['b_vocab'] = np.zeros(vocab_size)

        # Cast parameters to correct dtype
```

```
 74          for k, v in self.params.items():
 75              self.params[k] = v.astype(self.dtype)
 76
 77
 78      def loss(self, features, captions):
 79          """
 80          Compute training-time loss for the RNN. We input image features and
 81          ground-truth captions for those images, and use an RNN (or LSTM) to compute
 82          loss and gradients on all parameters.
 83
 84          Inputs:
 85          - features: Input image features, of shape (N, D)
 86          - captions: Ground-truth captions; an integer array of shape (N, T) where
 87            each element is in the range 0 <= y[i, t] < V
 88
 89          Returns a tuple of:
 90          - loss: Scalar loss
 91          - grads: Dictionary of gradients parallel to self.params
 92          """
 93          # Cut captions into two pieces: captions_in has everything but the last word
 94          # and will be input to the RNN; captions_out has everything but the first
 95          # word and this is what we will expect the RNN to generate. These are offset
 96          # by one relative to each other because the RNN should produce word (t+1)
 97          # after receiving word t. The first element of captions_in will be the START
 98          # token, and the first element of captions_out will be the first word.
 99          captions_in = captions[:, :-1]
100          captions_out = captions[:, 1:]
101
102          # You'll need this
103          mask = (captions_out != self._null)
104
105          # Weight and bias for the affine transform from image features to initial
106          # hidden state
107          W_proj, b_proj = self.params['W_proj'], self.params['b_proj']
108
109          # Word embedding matrix
110          W_embed = self.params['W_embed']
111
112          # Input-to-hidden, hidden-to-hidden, and biases for the RNN
113          Wx, Wh, b = self.params['Wx'], self.params['Wh'], self.params['b']
114
115          # Weight and bias for the hidden-to-vocab transformation.
116          W_vocab, b_vocab = self.params['W_vocab'], self.params['b_vocab']
117
118          loss, grads = 0.0, {}
119          ############################################################################
120          # TODO: Implement the forward and backward passes for the CaptioningRNN.   #
121          # In the forward pass you will need to do the following:                   #
122          # (1) Use an affine transformation to compute the initial hidden state     #
123          #     from the image features. This should produce an array of shape (N, H)#
124          # (2) Use a word embedding layer to transform the words in captions_in     #
125          #     from indices to vectors, giving an array of shape (N, T, W).         #
126          # (3) Use either a vanilla RNN or LSTM (depending on self.cell_type) to    #
127          #     process the sequence of input word vectors and produce hidden state  #
128          #     vectors for all timesteps, producing an array of shape (N, T, H).    #
129          # (4) Use a (temporal) affine transformation to compute scores over the    #
130          #     vocabulary at every timestep using the hidden states, giving an      #
131          #     array of shape (N, T, V).                                            #
132          # (5) Use (temporal) softmax to compute loss using captions_out, ignoring  #
133          #     the points where the output word is <NULL> using the mask above.     #
134          #                                                                          #
135          # In the backward pass you will need to compute the gradient of the loss   #
136          # with respect to all model parameters. Use the loss and grads variables   #
137          # defined above to store loss and gradients; grads[k] should give the      #
138          # gradients for self.params[k].                                            #
139          #                                                                          #
140          # Note also that you are allowed to make use of functions from layers.py   #
141          # in your implementation, if needed.                                       #
142          ############################################################################
143          # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
144
145          # Forward pass
146          initial_hidden, initial_cache = affine_forward(features, W_proj, b_proj)
147          captions_embedded, embedding_cache = word_embedding_forward(captions_in, W_embed)
148          if self.cell_type == "rnn":
149              rnn_output, rnn_cache = rnn_forward(captions_embedded, initial_hidden, Wx, Wh, b)
```

```
150         elif self.cell_type == 'lstm':
151             rnn_output, lstm_cache = lstm_forward(captions_embedded, initial_hidden, Wx, Wh, b)
152
153         scores, vocab_cache = temporal_affine_forward(rnn_output, W_vocab, b_vocab)
154         loss, dx = temporal_softmax_loss(scores, captions_out, mask)
155
156         # Backward pass
157         dscores, grads['W_vocab'], grads['b_vocab'] = temporal_affine_backward(dx, vocab_cache)
158
159         if self.cell_type == 'rnn':
160             dx, dh_initial, grads['Wx'], grads['Wh'], grads['b'] = rnn_backward(dscores, rnn_cache)
161         elif self.cell_type == 'lstm':
162             dx, dh_initial, grads['Wx'], grads['Wh'], grads['b'] = lstm_backward(dscores, lstm_cache)
163
164         grads['W_embed'] = word_embedding_backward(dx, embedding_cache)
165         dx_initial, grads['W_proj'], grads['b_proj'] = affine_backward(dh_initial, initial_cache)
166
167         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
168         ############################################################################
169         #                            END OF YOUR CODE                              #
170         ############################################################################
171
172         return loss, grads
173
174
175     def sample(self, features, max_length=30):
176         """
177         Run a test-time forward pass for the model, sampling captions for input
178         feature vectors.
179
180         At each timestep, we embed the current word, pass it and the previous hidden
181         state to the RNN to get the next hidden state, use the hidden state to get
182         scores for all vocab words, and choose the word with the highest score as
183         the next word. The initial hidden state is computed by applying an affine
184         transform to the input image features, and the initial word is the <START>
185         token.
186
187         For LSTMs you will also have to keep track of the cell state; in that case
188         the initial cell state should be zero.
189
190         Inputs:
191         - features: Array of input image features of shape (N, D).
192         - max_length: Maximum length T of generated captions.
193
194         Returns:
195         - captions: Array of shape (N, max_length) giving sampled captions,
196           where each element is an integer in the range [0, V). The first element
197           of captions should be the first sampled word, not the <START> token.
198         """
199         N = features.shape[0]
200         captions = self._null * np.ones((N, max_length), dtype=np.int32)
201
202         # Unpack parameters
203         W_proj, b_proj = self.params['W_proj'], self.params['b_proj']
204         W_embed = self.params['W_embed']
205         Wx, Wh, b = self.params['Wx'], self.params['Wh'], self.params['b']
206         W_vocab, b_vocab = self.params['W_vocab'], self.params['b_vocab']
207
208         ############################################################################
209         # TODO: Implement test-time sampling for the model. You will need to       #
210         # initialize the hidden state of the RNN by applying the learned affine    #
211         # transform to the input image features. The first word that you feed to   #
212         # the RNN should be the <START> token; its value is stored in the          #
213         # variable self._start. At each timestep you will need to do to:           #
214         # (1) Embed the previous word using the learned word embeddings            #
215         # (2) Make an RNN step using the previous hidden state and the embedded     #
216         #     current word to get the next hidden state.                           #
217         # (3) Apply the learned affine transformation to the next hidden state to  #
218         #     get scores for all words in the vocabulary                           #
219         # (4) Select the word with the highest score as the next word, writing it  #
220         #     (the word index) to the appropriate slot in the captions variable    #
221         #                                                                          #
222         # For simplicity, you do not need to stop generating after an <END> token  #
223         # is sampled, but you can if you want to.                                  #
224         #                                                                          #
225         # HINT: You will not be able to use the rnn_forward or lstm_forward        #
```

```python
          # functions; you'll need to call rnn_step_forward or lstm_step_forward in  #
          # a loop.                                                                   #
          #                                                                           #
          # NOTE: we are still working over minibatches in this function. Also if     #
          # you are using an LSTM, initialize the first cell state to zeros.          #
          #############################################################################
          # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

          hidden_state, hidden_cache = affine_forward(features, W_proj, b_proj)

          if self.cell_type == 'lstm':
              cell_state = np.zeros_like(hidden_state)

          word_embed, embedding_cache = word_embedding_forward(self._start, W_embed)

          for i in range(max_length):
              if self.cell_type == 'rnn':
                  hidden_state, rnn_cache = rnn_step_forward(word_embed, hidden_state, Wx, Wh, b)
              elif self.cell_type == 'lstm':
                  hidden_state, cell_state, rnn_cache = lstm_step_forward(word_embed, hidden_state,
      cell_state, Wx, Wh, b)

              scores, affine_cache = affine_forward(hidden_state, W_vocab, b_vocab)

              captions[:,i] = np.argmax(scores, axis=1)

              word_embed, _ = word_embedding_forward(captions[:, i], W_embed)

          # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
          #############################################################################
          #                             END OF YOUR CODE                              #
          #############################################################################
          return captions
```

## 2 rnn_layers.py

```python
from __future__ import print_function, division
from builtins import range
import numpy as np


"""
This file defines layer types that are commonly used for recurrent neural
networks.
"""


def rnn_step_forward(x, prev_h, Wx, Wh, b):
    """
    Run the forward pass for a single timestep of a vanilla RNN that uses a tanh
    activation function.

    The input data has dimension D, the hidden state has dimension H, and we use
    a minibatch size of N.

    Inputs:
    - x: Input data for this timestep, of shape (N, D).
    - prev_h: Hidden state from previous timestep, of shape (N, H)
    - Wx: Weight matrix for input-to-hidden connections, of shape (D, H)
    - Wh: Weight matrix for hidden-to-hidden connections, of shape (H, H)
    - b: Biases of shape (H,)

    Returns a tuple of:
    - next_h: Next hidden state, of shape (N, H)
    - cache: Tuple of values needed for the backward pass.
    """
    next_h, cache = None, None
    ##########################################################################
    # TODO: Implement a single forward step for the vanilla RNN. Store the next  #
    # hidden state and any values you need for the backward pass in the next_h   #
    # and cache variables respectively.                                          #
    ##########################################################################
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    next_h = np.tanh(np.dot(x, Wx) + np.dot(prev_h, Wh) + b)
    cache = (x, prev_h, Wx, Wh, b, next_h)

    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
    ##########################################################################
    #                          END OF YOUR CODE                              #
    ##########################################################################
    return next_h, cache


def rnn_step_backward(dnext_h, cache):
    """
    Backward pass for a single timestep of a vanilla RNN.

    Inputs:
    - dnext_h: Gradient of loss with respect to next hidden state, of shape (N, H)
    - cache: Cache object from the forward pass

    Returns a tuple of:
    - dx: Gradients of input data, of shape (N, D)
    - dprev_h: Gradients of previous hidden state, of shape (N, H)
    - dWx: Gradients of input-to-hidden weights, of shape (D, H)
    - dWh: Gradients of hidden-to-hidden weights, of shape (H, H)
    - db: Gradients of bias vector, of shape (H,)
    """
    dx, dprev_h, dWx, dWh, db = None, None, None, None, None
    ##########################################################################
    # TODO: Implement the backward pass for a single step of a vanilla RNN.      #
    #                                                                            #
    # HINT: For the tanh function, you can compute the local derivative in terms #
    # of the output value from tanh.                                             #
    ##########################################################################
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    x, prev_h, Wx, Wh, b, next_h = cache
```

```
74        dhraw = (1 − next_h ∗ next_h) ∗ dnext_h
75        db = np.sum(dhraw, axis = 0)
76        dx = np.dot(dhraw, Wx.T)
77        dWx = np.dot(x.T, dhraw)
78        dprev_h = np.dot(dhraw, Wh.T)
79        dWh = np.dot(prev_h.T, dhraw)
80
81
82
83        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
84        ##############################################################################
85        #                            END OF YOUR CODE                                #
86        ##############################################################################
87        return dx, dprev_h, dWx, dWh, db
88
89
90  def rnn_forward(x, h0, Wx, Wh, b):
91        """
92        Run a vanilla RNN forward on an entire sequence of data. We assume an input
93        sequence composed of T vectors, each of dimension D. The RNN uses a hidden
94        size of H, and we work over a minibatch containing N sequences. After running
95        the RNN forward, we return the hidden states for all timesteps.
96
97        Inputs:
98        − x: Input data for the entire timeseries, of shape (N, T, D).
99        − h0: Initial hidden state, of shape (N, H)
100       − Wx: Weight matrix for input−to−hidden connections, of shape (D, H)
101       − Wh: Weight matrix for hidden−to−hidden connections, of shape (H, H)
102       − b: Biases of shape (H,)
103
104       Returns a tuple of:
105       − h: Hidden states for the entire timeseries, of shape (N, T, H).
106       − cache: Values needed in the backward pass
107       """
108       h, cache = None, None
109       ##############################################################################
110       # TODO: Implement forward pass for a vanilla RNN running on a sequence of     #
111       # input data. You should use the rnn_step_forward function that you defined   #
112       # above. You can use a for loop to help compute the forward pass.             #
113       ##############################################################################
114       # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
115
116       N, T, D= x.shape
117       H = h0.shape[1]
118
119       h = np.zeros([N, T, H])
120       cache = []
121
122
123       for i in range(T):
124           next_h, cache_i = rnn_step_forward(x[:, i, :], h0, Wx, Wh, b)
125           h0 = next_h
126           h[:, i, :] = next_h
127           cache.append(cache_i)
128
129       # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
130       ##############################################################################
131       #                            END OF YOUR CODE                                #
132       ##############################################################################
133       return h, cache
134
135
136 def rnn_backward(dh, cache):
137       """
138       Compute the backward pass for a vanilla RNN over an entire sequence of data.
139
140       Inputs:
141       − dh: Upstream gradients of all hidden states, of shape (N, T, H).
142
143       NOTE: 'dh' contains the upstream gradients produced by the
144       individual loss functions at each timestep, *not* the gradients
145       being passed between timesteps (which you'll have to compute yourself
146       by calling rnn_step_backward in a loop).
147
148       Returns a tuple of:
149       − dx: Gradient of inputs, of shape (N, T, D)
```

```python
150         - dh0: Gradient of initial hidden state, of shape (N, H)
151         - dWx: Gradient of input-to-hidden weights, of shape (D, H)
152         - dWh: Gradient of hidden-to-hidden weights, of shape (H, H)
153         - db: Gradient of biases, of shape (H,)
154         """
155         dx, dh0, dWx, dWh, db = None, None, None, None, None
156         ###########################################################################
157         # TODO: Implement the backward pass for a vanilla RNN running an entire    #
158         # sequence of data. You should use the rnn_step_backward function that you #
159         # defined above. You can use a for loop to help compute the backward pass. #
160         ###########################################################################
161         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
162
163         x, prev_h, Wx, Wh, b, next_h = cache[0]
164         N, T, H = dh.shape
165         D = Wx.shape[0]
166
167         dx = np.zeros([N, T, D])
168         dWx = np.zeros_like(Wx)
169         dWh = np.zeros_like(Wh)
170         db = np.zeros_like(b)
171         dprev_h = np.zeros_like(prev_h)
172
173         dh0 = np.zeros([N, H])
174
175         for i in reversed(range(T)):
176             dh_i = dprev_h + dh[:,i,:]
177             dx[:, i, :], dprev_h, dWx_i, dWh_i, db_i = rnn_step_backward(dh_i, cache[i])
178             dWx += dWx_i
179             dWh += dWh_i
180             db += db_i
181
182         dh0 = dprev_h
183
184
185
186         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
187         ###########################################################################
188         #                             END OF YOUR CODE                            #
189         ###########################################################################
190         return dx, dh0, dWx, dWh, db
191
192
193 def word_embedding_forward(x, W):
194     """
195     Forward pass for word embeddings. We operate on minibatches of size N where
196     each sequence has length T. We assume a vocabulary of V words, assigning each
197     word to a vector of dimension D.
198
199     Inputs:
200     - x: Integer array of shape (N, T) giving indices of words. Each element idx
201       of x muxt be in the range 0 <= idx < V.
202     - W: Weight matrix of shape (V, D) giving word vectors for all words.
203
204     Returns a tuple of:
205     - out: Array of shape (N, T, D) giving word vectors for all input words.
206     - cache: Values needed for the backward pass
207     """
208     out, cache = None, None
209     ###########################################################################
210     # TODO: Implement the forward pass for word embeddings.                   #
211     #                                                                         #
212     # HINT: This can be done in one line using NumPy's array indexing.        #
213     ###########################################################################
214     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
215
216     # out are values chosen from W with corresponding indices provided in x
217     out = W[x, :]
218     cache = x, W
219
220     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
221     ###########################################################################
222     #                             END OF YOUR CODE                            #
223     ###########################################################################
224     return out, cache
225
```

```
226
227  def word_embedding_backward(dout, cache):
228      """
229      Backward pass for word embeddings. We cannot back-propagate into the words
230      since they are integers, so we only return gradient for the word embedding
231      matrix.
232
233      HINT: Look up the function np.add.at
234
235      Inputs:
236      - dout: Upstream gradients of shape (N, T, D)
237      - cache: Values from the forward pass
238
239      Returns:
240      - dW: Gradient of word embedding matrix, of shape (V, D).
241      """
242      dW = None
243      ###########################################################################
244      # TODO: Implement the backward pass for word embeddings.                  #
245      #                                                                         #
246      # Note that words can appear more than once in a sequence.                #
247      # HINT: Look up the function np.add.at                                    #
248      ###########################################################################
249      # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
250
251      x, W = cache
252      V = W.shape[0]
253      dW = np.zeros_like(W)
254
255      # Add dout of the corresponding indices (x) to dW
256      np.add.at(dW, x, dout)
257
258      # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
259      ###########################################################################
260      #                          END OF YOUR CODE                               #
261      ###########################################################################
262      return dW
263
264
265  def sigmoid(x):
266      """
267      A numerically stable version of the logistic sigmoid function.
268      """
269      pos_mask = (x >= 0)
270      neg_mask = (x < 0)
271      z = np.zeros_like(x)
272      z[pos_mask] = np.exp(-x[pos_mask])
273      z[neg_mask] = np.exp(x[neg_mask])
274      top = np.ones_like(x)
275      top[neg_mask] = z[neg_mask]
276      return top / (1 + z)
277
278
279  def lstm_step_forward(x, prev_h, prev_c, Wx, Wh, b):
280      """
281      Forward pass for a single timestep of an LSTM.
282
283      The input data has dimension D, the hidden state has dimension H, and we use
284      a minibatch size of N.
285
286      Note that a sigmoid() function has already been provided for you in this file.
287
288      Inputs:
289      - x: Input data, of shape (N, D)
290      - prev_h: Previous hidden state, of shape (N, H)
291      - prev_c: previous cell state, of shape (N, H)
292      - Wx: Input-to-hidden weights, of shape (D, 4H)
293      - Wh: Hidden-to-hidden weights, of shape (H, 4H)
294      - b: Biases, of shape (4H,)
295
296      Returns a tuple of:
297      - next_h: Next hidden state, of shape (N, H)
298      - next_c: Next cell state, of shape (N, H)
299      - cache: Tuple of values needed for backward pass.
300      """
301      next_h, next_c, cache = None, None, None
```

```
302    ######################################################################################
303    # TODO: Implement the forward pass for a single timestep of an LSTM.              #
304    # You may want to use the numerically stable sigmoid implementation above.        #
305    ######################################################################################
306    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
307
308    A = x.dot(Wx) + prev_h.dot(Wh) + b
309
310    H = prev_h.shape[1]
311
312    ai = A[:, 0: H]
313    af = A[:, H: 2*H]
314    ao = A[:, 2*H: 3*H]
315    ag = A[:, 3*H: 4*H]
316
317    i = sigmoid(ai)
318    f = sigmoid(af)
319    o = sigmoid(ao)
320    g = np.tanh(ag)
321
322    next_c = f * prev_c + i * g
323    next_h = o * np.tanh(next_c)
324
325    cache = (x, prev_h, prev_c, Wx, Wh, b, i, f, o, g, next_c, next_h)
326
327    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
328    ######################################################################################
329    #                            END OF YOUR CODE                                      #
330    ######################################################################################
331
332    return next_h, next_c, cache
333
334
335 def lstm_step_backward(dnext_h, dnext_c, cache):
336     """
337     Backward pass for a single timestep of an LSTM.
338
339     Inputs:
340     - dnext_h: Gradients of next hidden state, of shape (N, H)
341     - dnext_c: Gradients of next cell state, of shape (N, H)
342     - cache: Values from the forward pass
343
344     Returns a tuple of:
345     - dx: Gradient of input data, of shape (N, D)
346     - dprev_h: Gradient of previous hidden state, of shape (N, H)
347     - dprev_c: Gradient of previous cell state, of shape (N, H)
348     - dWx: Gradient of input-to-hidden weights, of shape (D, 4H)
349     - dWh: Gradient of hidden-to-hidden weights, of shape (H, 4H)
350     - db: Gradient of biases, of shape (4H,)
351     """
352     dx, dprev_h, dprev_c, dWx, dWh, db = None, None, None, None, None, None
353     ######################################################################################
354     # TODO: Implement the backward pass for a single timestep of an LSTM.              #
355     #                                                                                  #
356     # HINT: For sigmoid and tanh you can compute local derivatives in terms of         #
357     # the output value from the nonlinearity.                                          #
358     ######################################################################################
359     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
360
361     x, prev_h, prev_c, Wx, Wh, b, i, f, o, g, next_c, next_h = cache
362     do = np.tanh(next_c) * dnext_h
363     dnext_c += o * (1 - np.square(np.tanh(next_c))) * dnext_h
364     df = prev_c * dnext_c
365     dprev_c = f * dnext_c
366     di = g * dnext_c
367     dg = i * dnext_c
368
369
370     N, H = dnext_h.shape
371     dA = np.zeros((N, 4*H))
372     dA[:,0:H] = di * i *(1 - i)
373     dA[:,H:2*H] = df * f * (1 - f)
374     dA[:,2*H:3*H] = do * o * (1 - o)
375     dA[:,3*H:] = dg * (1 - np.square(g))
376
377     dx = np.dot(dA, Wx.T)
```

```
378        dWx = np.dot(x.T, dA)
379        dprev_h = np.dot(dA, Wh.T)
380        dWh = np.dot(prev_h.T, dA)
381        db = np.sum(dA, axis=0)



385        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
386        ##############################################################################
387        #                           END OF YOUR CODE                                 #
388        ##############################################################################

390        return dx, dprev_h, dprev_c, dWx, dWh, db


393    def lstm_forward(x, h0, Wx, Wh, b):
394        """
395        Forward pass for an LSTM over an entire sequence of data. We assume an input
396        sequence composed of T vectors, each of dimension D. The LSTM uses a hidden
397        size of H, and we work over a minibatch containing N sequences. After running
398        the LSTM forward, we return the hidden states for all timesteps.
399
400        Note that the initial cell state is passed as input, but the initial cell
401        state is set to zero. Also note that the cell state is not returned; it is
402        an internal variable to the LSTM and is not accessed from outside.
403
404        Inputs:
405        - x: Input data of shape (N, T, D)
406        - h0: Initial hidden state of shape (N, H)
407        - Wx: Weights for input-to-hidden connections, of shape (D, 4H)
408        - Wh: Weights for hidden-to-hidden connections, of shape (H, 4H)
409        - b: Biases of shape (4H,)
410
411        Returns a tuple of:
412        - h: Hidden states for all timesteps of all sequences, of shape (N, T, H)
413        - cache: Values needed for the backward pass.
414        """
415        h, cache = None, None
416        ##############################################################################
417        # TODO: Implement the forward pass for an LSTM over an entire timeseries.    #
418        # You should use the lstm_step_forward function that you just defined.       #
419        ##############################################################################
420        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

422        N, T, D= x.shape
423        H = h0.shape[1]

425        prev_h = h0
426        prev_c = np.zeros_like(h0)

428        h = np.zeros([N, T, H])
429        cache = []


432        for i in range(T):
433            next_h, next_c, cache_i = lstm_step_forward(x[:, i, :], prev_h, prev_c, Wx, Wh, b)
434            prev_h = next_h
435            prev_c = next_c
436            h[:, i, :] = next_h
437            cache.append(cache_i)

439        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
440        ##############################################################################
441        #                           END OF YOUR CODE                                 #
442        ##############################################################################

444        return h, cache


447    def lstm_backward(dh, cache):
448        """
449        Backward pass for an LSTM over an entire sequence of data.]
450
451        Inputs:
452        - dh: Upstream gradients of hidden states, of shape (N, T, H)
453        - cache: Values from the forward pass
```

```
454
455        Returns a tuple of:
456        - dx: Gradient of input data of shape (N, T, D)
457        - dh0: Gradient of initial hidden state of shape (N, H)
458        - dWx: Gradient of input-to-hidden weight matrix of shape (D, 4H)
459        - dWh: Gradient of hidden-to-hidden weight matrix of shape (H, 4H)
460        - db: Gradient of biases, of shape (4H,)
461        """
462        dx, dh0, dWx, dWh, db = None, None, None, None, None
463        ###########################################################################
464        # TODO: Implement the backward pass for an LSTM over an entire timeseries. #
465        # You should use the lstm_step_backward function that you just defined.    #
466        ###########################################################################
467        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
468
469        x, prev_h, prev_c, Wx, Wh, b, i, f, o, g, next_c, next_h = cache[0]
470        N, T, H = dh.shape
471        D = Wx.shape[0]
472
473        dx = np.zeros([N, T, D])
474        dWx = np.zeros_like(Wx)
475        dWh = np.zeros_like(Wh)
476        db = np.zeros_like(b)
477        dprev_h = np.zeros_like(prev_h)
478        dprev_c = np.zeros_like(prev_c)
479
480        dh0 = np.zeros([N, H])
481
482        for i in reversed(range(T)):
483            dh_i = dprev_h + dh[:,i,:]
484
485            dx[:,i,:], dprev_h, dprev_c, dWx_i, dWh_i, db_i = lstm_step_backward(dh_i, dprev_c, cache[i])
486
487            db += db_i
488            dWh += dWh_i
489            dWx += dWx_i
490
491
492        dh0 = dprev_h
493
494
495
496        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
497        ###########################################################################
498        #                            END OF YOUR CODE                             #
499        ###########################################################################
500
501        return dx, dh0, dWx, dWh, db
502
503
504 def temporal_affine_forward(x, w, b):
505        """
506        Forward pass for a temporal affine layer. The input is a set of D-dimensional
507        vectors arranged into a minibatch of N timeseries, each of length T. We use
508        an affine function to transform each of those vectors into a new vector of
509        dimension M.
510
511        Inputs:
512        - x: Input data of shape (N, T, D)
513        - w: Weights of shape (D, M)
514        - b: Biases of shape (M,)
515
516        Returns a tuple of:
517        - out: Output data of shape (N, T, M)
518        - cache: Values needed for the backward pass
519        """
520        N, T, D = x.shape
521        M = b.shape[0]
522        out = x.reshape(N * T, D).dot(w).reshape(N, T, M) + b
523        cache = x, w, b, out
524        return out, cache
525
526
527 def temporal_affine_backward(dout, cache):
528        """
529        Backward pass for temporal affine layer.
```

```
530
531         Input:
532         - dout: Upstream gradients of shape (N, T, M)
533         - cache: Values from forward pass
534
535         Returns a tuple of:
536         - dx: Gradient of input, of shape (N, T, D)
537         - dw: Gradient of weights, of shape (D, M)
538         - db: Gradient of biases, of shape (M,)
539         """
540         x, w, b, out = cache
541         N, T, D = x.shape
542         M = b.shape[0]
543
544         dx = dout.reshape(N * T, M).dot(w.T).reshape(N, T, D)
545         dw = dout.reshape(N * T, M).T.dot(x.reshape(N * T, D)).T
546         db = dout.sum(axis=(0, 1))
547
548         return dx, dw, db
549
550
551  def temporal_softmax_loss(x, y, mask, verbose=False):
552         """
553         A temporal version of softmax loss for use in RNNs. We assume that we are
554         making predictions over a vocabulary of size V for each timestep of a
555         timeseries of length T, over a minibatch of size N. The input x gives scores
556         for all vocabulary elements at all timesteps, and y gives the indices of the
557         ground-truth element at each timestep. We use a cross-entropy loss at each
558         timestep, summing the loss over all timesteps and averaging across the
559         minibatch.
560
561         As an additional complication, we may want to ignore the model output at some
562         timesteps, since sequences of different length may have been combined into a
563         minibatch and padded with NULL tokens. The optional mask argument tells us
564         which elements should contribute to the loss.
565
566         Inputs:
567         - x: Input scores, of shape (N, T, V)
568         - y: Ground-truth indices, of shape (N, T) where each element is in the range
569             0 <= y[i, t] < V
570         - mask: Boolean array of shape (N, T) where mask[i, t] tells whether or not
571           the scores at x[i, t] should contribute to the loss.
572
573         Returns a tuple of:
574         - loss: Scalar giving loss
575         - dx: Gradient of loss with respect to scores x.
576         """
577
578         N, T, V = x.shape
579
580         x_flat = x.reshape(N * T, V)
581         y_flat = y.reshape(N * T)
582         mask_flat = mask.reshape(N * T)
583
584         probs = np.exp(x_flat - np.max(x_flat, axis=1, keepdims=True))
585         probs /= np.sum(probs, axis=1, keepdims=True)
586         loss = -np.sum(mask_flat * np.log(probs[np.arange(N * T), y_flat])) / N
587         dx_flat = probs.copy()
588         dx_flat[np.arange(N * T), y_flat] -= 1
589         dx_flat /= N
590         dx_flat *= mask_flat[:, None]
591
592         if verbose: print('dx_flat: ', dx_flat.shape)
593
594         dx = dx_flat.reshape(N, T, V)
595
596         return loss, dx
```