

RNN_Captioning

December 4, 2020

1 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.

1.1 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:

```
[ ]: !pip install h5py
```

```
[1]: # As usual, a bit of setup
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient,
    eval_numerical_gradient_array
from cs231n.rnn_layers import *
from cs231n.captioning_solver import CaptioningSolver
from cs231n.classifiers.rnn import CaptioningRNN
from cs231n.coco_utils import load_coco_data, sample_coco_minibatch,
    decode_captions
from cs231n.image_utils import image_from_url

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
    ↵autoreload-of-modules-in-ipython
```

```
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

2 Microsoft COCO

For this exercise we will use the 2014 release of the [Microsoft COCO dataset](#) which has become the standard testbed for image captioning. The dataset consists of 80,000 training images and 40,000 validation images, each annotated with 5 captions written by workers on Amazon Mechanical Turk.

You should have already downloaded the data by changing to the `cs231n/datasets` directory and running the script `get_assignment3_data.sh`. If you haven't yet done so, run that script now. Warning: the COCO data download is ~1GB.

We have preprocessed the data and extracted features for you already. For all images we have extracted features from the fc7 layer of the VGG-16 network pretrained on ImageNet; these features are stored in the files `train2014_vgg16_fc7.h5` and `val2014_vgg16_fc7.h5` respectively. To cut down on processing time and memory requirements, we have reduced the dimensionality of the features from 4096 to 512; these features can be found in the files `train2014_vgg16_fc7_pca.h5` and `val2014_vgg16_fc7_pca.h5`.

The raw images take up a lot of space (nearly 20GB) so we have not included them in the download. However all images are taken from Flickr, and URLs of the training and validation images are stored in the files `train2014_urls.txt` and `val2014_urls.txt` respectively. This allows you to download images on the fly for visualization. Since images are downloaded on-the-fly, **you must be connected to the internet to view images**.

Dealing with strings is inefficient, so we will work with an encoded version of the captions. Each word is assigned an integer ID, allowing us to represent a caption by a sequence of integers. The mapping between integer IDs and words is in the file `coco2014_vocab.json`, and you can use the function `decode_captions` from the file `cs231n/coco_utils.py` to convert numpy arrays of integer IDs back into strings.

There are a couple special tokens that we add to the vocabulary. We prepend a special `<START>` token and append an `<END>` token to the beginning and end of each caption respectively. Rare words are replaced with a special `<UNK>` token (for “unknown”). In addition, since we want to train with minibatches containing captions of different lengths, we pad short captions with a special `<NULL>` token after the `<END>` token and don’t compute loss or gradient for `<NULL>` tokens. Since they are a bit of a pain, we have taken care of all implementation details around special tokens for you.

You can load all of the MS-COCO data (captions, features, URLs, and vocabulary) using the `load_coco_data` function from the file `cs231n/coco_utils.py`. Run the following cell to do so:

```
[2]: # 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))

base dir
/Users/luciabouza/FING/DLVis/assignment3_jupyter/cs231n/datasets/coco_captioning
trainCaptions <class 'numpy.ndarray'> (400135, 17) int32
trainImageIdxs <class 'numpy.ndarray'> (400135,) int32
valCaptions <class 'numpy.ndarray'> (195954, 17) int32
valImageIdxs <class 'numpy.ndarray'> (195954,) int32
trainFeatures <class 'numpy.ndarray'> (82783, 512) float32
valFeatures <class 'numpy.ndarray'> (40504, 512) float32
idx_to_word <class 'list'> 1004
word_to_idx <class 'dict'> 1004
trainUrls <class 'numpy.ndarray'> (82783,) <U63
valUrls <class 'numpy.ndarray'> (40504,) <U63

```

2.1 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**.

```
[49]: # 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 dark <UNK> woman holding a knife to her <UNK> <END>



<START> a photo taken from a field looking at a train going by <END>



<START> three surfers enjoying a large ocean wave together <END>



3 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`.

4 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.

```
[6]: 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

5 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.

```
[11]: 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:  2.319932372313319e-10
dprev_h error:  2.6828355645784327e-10
dWx error:  8.820244454238703e-10
dWh error:  4.703287554560559e-10
db error:  1.5956895526227225e-11

```

6 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.

```

[15]: 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]]])
print('h error: ', rel_error(expected_h, h))

```

h error: 7.728466151011529e-08

7 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.

```
[26]: 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, 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: 1.5402322184213243e-09
dh0 error: 3.3824326261334578e-09
dWx error: 7.238350796069372e-09
dWh error: 1.3157659173166636e-07
db error: 1.5353591509855146e-10

```

8 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.

[27]: 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
```

9 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.

```
[32]: 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

10 Temporal Affine layer

At every timestep we use an affine function to transform the RNN hidden vector at that timestep into scores for each word in the vocabulary. Because this is very similar to the affine layer that you implemented in assignment 2, we have provided this function for you in the `temporal_affine_forward` and `temporal_affine_backward` functions in the file `cs231n/rnn_layers.py`. Run the following to perform numeric gradient checking on the implementation. You should see errors on the order of e-9 or less.

```
[33]: np.random.seed(231)

# Gradient check for temporal affine layer
N, T, D, M = 2, 3, 4, 5
x = np.random.randn(N, T, D)
w = np.random.randn(D, M)
b = np.random.randn(M)

out, cache = temporal_affine_forward(x, w, b)

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

fx = lambda x: temporal_affine_forward(x, w, b)[0]
fw = lambda w: temporal_affine_forward(x, w, b)[0]
fb = lambda b: temporal_affine_forward(x, w, b)[0]
```

```

dx_num = eval_numerical_gradient_array(fx, x, dout)
dw_num = eval_numerical_gradient_array(fw, w, dout)
db_num = eval_numerical_gradient_array(fb, b, dout)

dx, dw, db = temporal_affine_backward(dout, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))

```

```

dx error:  2.9215945034030545e-10
dw error:  1.5772088618663602e-10
db error:  3.252200556967514e-11

```

11 Temporal Softmax loss

In an RNN language model, at every timestep we produce a score for each word in the vocabulary. We know the ground-truth word at each timestep, so we use a softmax loss function to compute loss and gradient at each timestep. We sum the losses over time and average them over the minibatch.

However there is one wrinkle: since we operate over minibatches and different captions may have different lengths, we append <NULL> tokens to the end of each caption so they all have the same length. We don't want these <NULL> tokens to count toward the loss or gradient, so in addition to scores and ground-truth labels our loss function also accepts a `mask` array that tells it which elements of the scores count towards the loss.

Since this is very similar to the softmax loss function you implemented in assignment 1, we have implemented this loss function for you; look at the `temporal_softmax_loss` function in the file `cs231n/rnn_layers.py`.

Run the following cell to sanity check the loss and perform numeric gradient checking on the function. You should see an error for `dx` on the order of e-7 or less.

```
[34]: # Sanity check for temporal softmax loss
from cs231n.rnn_layers import temporal_softmax_loss

N, T, V = 100, 1, 10

def check_loss(N, T, V):
    x = 0.001 * np.random.randn(N, T, V)
    y = np.random.randint(V, size=(N, T))
    mask = np.random.rand(N, T) <= p
    print(temporal_softmax_loss(x, y, mask)[0])

check_loss(100, 1, 10, 1.0)  # Should be about 2.3
check_loss(100, 10, 10, 1.0) # Should be about 23
check_loss(5000, 10, 10, 0.1) # Should be within 2.2-2.4

# Gradient check for temporal softmax loss
```

```

N, T, V = 7, 8, 9

x = np.random.randn(N, T, V)
y = np.random.randint(V, size=(N, T))
mask = (np.random.rand(N, T) > 0.5)

loss, dx = temporal_softmax_loss(x, y, mask, verbose=False)

dx_num = eval_numerical_gradient(lambda x: temporal_softmax_loss(x, y, mask)[0], x, verbose=False)

print('dx error: ', rel_error(dx, dx_num))

```

```

2.3027781774290146
23.025985953127226
2.2643611790293394
dx error: 2.583585303524283e-08

```

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.

```

[36]: 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: 9.832355910027387
expected loss: 9.83235591003
difference: 2.6130209107577684e-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.

```

[65]: 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: 2.331071e-09
W_proj relative error: 9.974426e-09
W_vocab relative error: 4.274378e-09
Wh relative error: 5.557955e-09
Wx relative error: 7.725620e-07

```

```
b relative error: 8.001353e-10
b_proj relative error: 6.260039e-09
b_vocab relative error: 1.690334e-09
```

13 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.

```
[39]: 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,
)

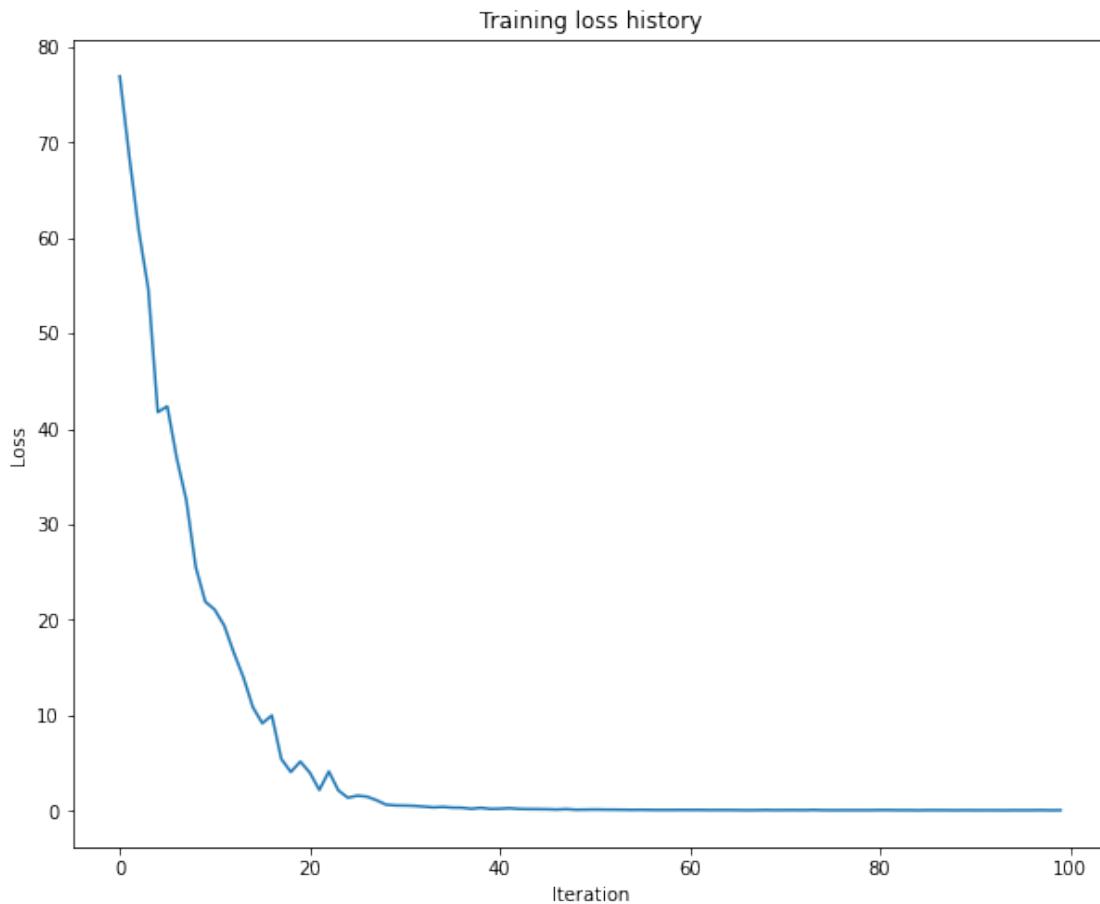
small_rnn_solver = CaptioningSolver(small_rnn_model, small_data,
    update_rule='adam',
    num_epochs=50,
    batch_size=25,
    optim_config={
        'learning_rate': 5e-3,
    },
    lr_decay=0.95,
    verbose=True, print_every=10,
)

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()

base dir
/Users/luciabouza/FING/DLVis/assignment3_jupyter/cs231n/datasets/coco_captioning
```

```
(Iteration 1 / 100) loss: 76.913486
(Iteration 11 / 100) loss: 21.063153
(Iteration 21 / 100) loss: 4.016260
(Iteration 31 / 100) loss: 0.567127
(Iteration 41 / 100) loss: 0.239437
(Iteration 51 / 100) loss: 0.162021
(Iteration 61 / 100) loss: 0.111542
(Iteration 71 / 100) loss: 0.097580
(Iteration 81 / 100) loss: 0.099095
(Iteration 91 / 100) loss: 0.073979
```



Print final training loss. You should see a final loss of less than 0.1.

```
[40]: print('Final loss: ', small_rnn_solver.loss_history[-1])
```

```
Final loss: 0.08208770288494388
```

14 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.

```
[63]: 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 large truck parked on the cement in front of a building <END>
GT:<START> a large truck parked on the cement in front of a building <END>



train

a <UNK> and a pair of shoes <UNK> together outside <END>
GT:<START> a <UNK> and a pair of shoes <UNK> together outside <END>



val

kids with a <END>

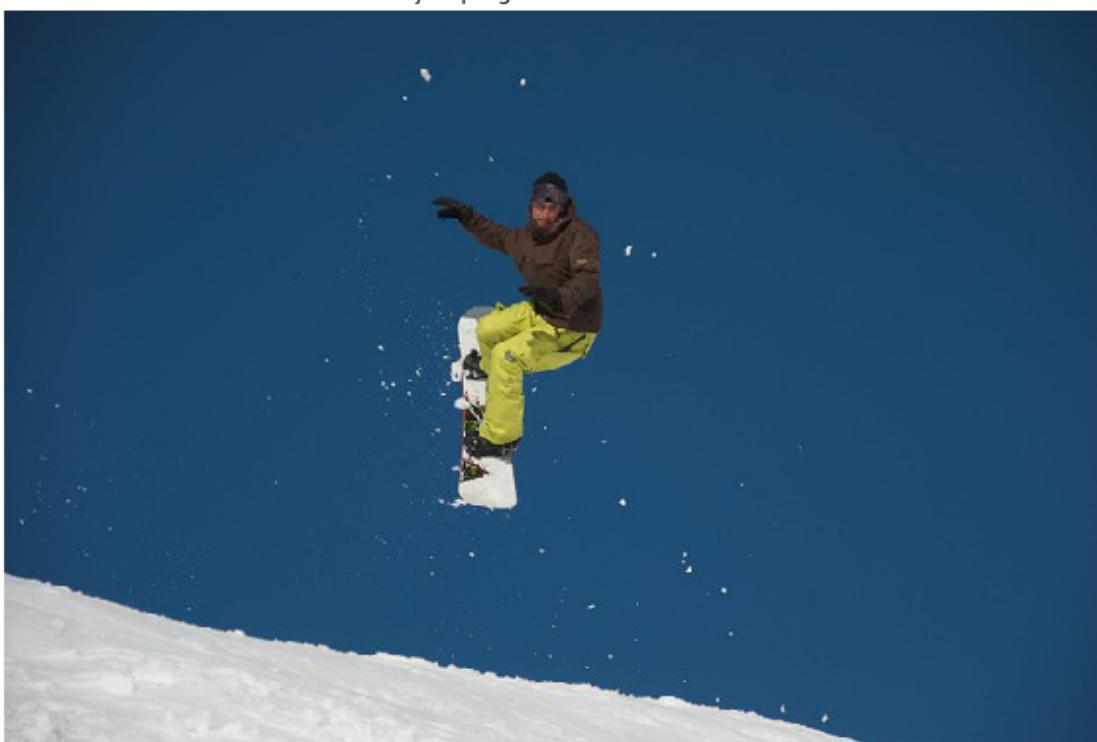
GT:<START> a kitchen with pots hanging over the stove <END>



val

half of <UNK> <UNK> in in clouds clouds <END>

GT:<START> someone jumping in the air on their snowboard <END>



15 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 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:

Al realizar un modelo de predicción por carácter, el vocabulario se reduce drásticamente. Mientras que si utilizo un modelo de palabras, el vocabulario superará probablemente las 10.000 palabras. En un modelo por carácter, será órdenes menor. Tener un vocabulario menor implica que las matrices de pesos W_{embed} y W_{vocab} serán menores, así como los vectores de bias. En un modelo de predicción por carácter tendremos menos parámetros para aprender. Si tenemos poca capacidad de cómputo o pocos datos, quizás un modelo por carácter pueda funcionar mejor.

Sin embargo, teniendo datos de entrenamiento y capacidad de cómputo, un modelo de predicción por palabra parece ser una mejor opción. Las RNN no pueden retener tanto contexto, digamos que no tienen tanta “memoria”, por lo que seguramente con un modelo por carácter, pierda contexto de la oración que está generando más fácilmente, y se generen oraciones sin sentido.

[]:

LSTM_Captioning

December 4, 2020

1 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.

```
[1]: # As usual, a bit of setup
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient, u
→eval_numerical_gradient_array
from cs231n.rnn_layers import *
from cs231n.captioning_solver import CaptioningSolver
from cs231n.classifiers.rnn import CaptioningRNN
from cs231n.coco_utils import load_coco_data, sample_coco_minibatch, u
→decode_captions
from cs231n.image_utils import image_from_url

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
→autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

2 Load MS-COCO data

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

```
[2]: # 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))
```

```
base dir
/Users/luciabouza/FING/DLVis/assignment3_jupyter/cs231n/datasets/coco_captioning
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
```

3 LSTM

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

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

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

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

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

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

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

where \odot is the elementwise product of vectors.

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

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

4 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.

```
[3]: 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:  5.7054131185818695e-09
next_c error:  5.8143123088804145e-09
```

5 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.

```
[7]: 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: 6.141307149471403e-10
dh error: 3.0914746081903265e-10
dc error: 1.5221723979041107e-10
dWx error: 1.6933647290937225e-09
dWh error: 4.80624861072581e-08
db error: 1.734923562619879e-10

```

6 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.

```

[8]: 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

```

7 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).

```
[10]: 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: 6.186242015576938e-09
dh0 error: 6.79192684354467e-09
dWx error: 3.301449189224954e-09
dWh error: 1.5077886632869585e-06
db error: 7.214760285323416e-10
```

8 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:

Los valores i , f y o siguen la idea de compuertas, que dejan o no dejar pasar información. En este caso las compuertas no son binarias, sino que tienen la posibilidad de dejar pasar una porción de información, por lo que se requiere que la función que las modele sea probabilista. La función sigmoide es probabilista, mientras que la función ReLu no lo es, y es por eso que habitualmente se utiliza la función sigmoide.

la función de cada compuerta es la siguiente: - La compuerta f nos indica cuánto de la información almacenada en la celda se mantiene o se borra. - La compuerta i nos indica cuánto de la información que ingresa se almacenará en la celda. - La compuerta o nos indica cuánto de la información almacenada en la celda se devolverá como salida.

9 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.

```
[11]: 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: 9.824459354432264
expected loss: 9.82445935443
difference: 2.2648549702353193e-12
```

10 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.

```
[12]: 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,
)

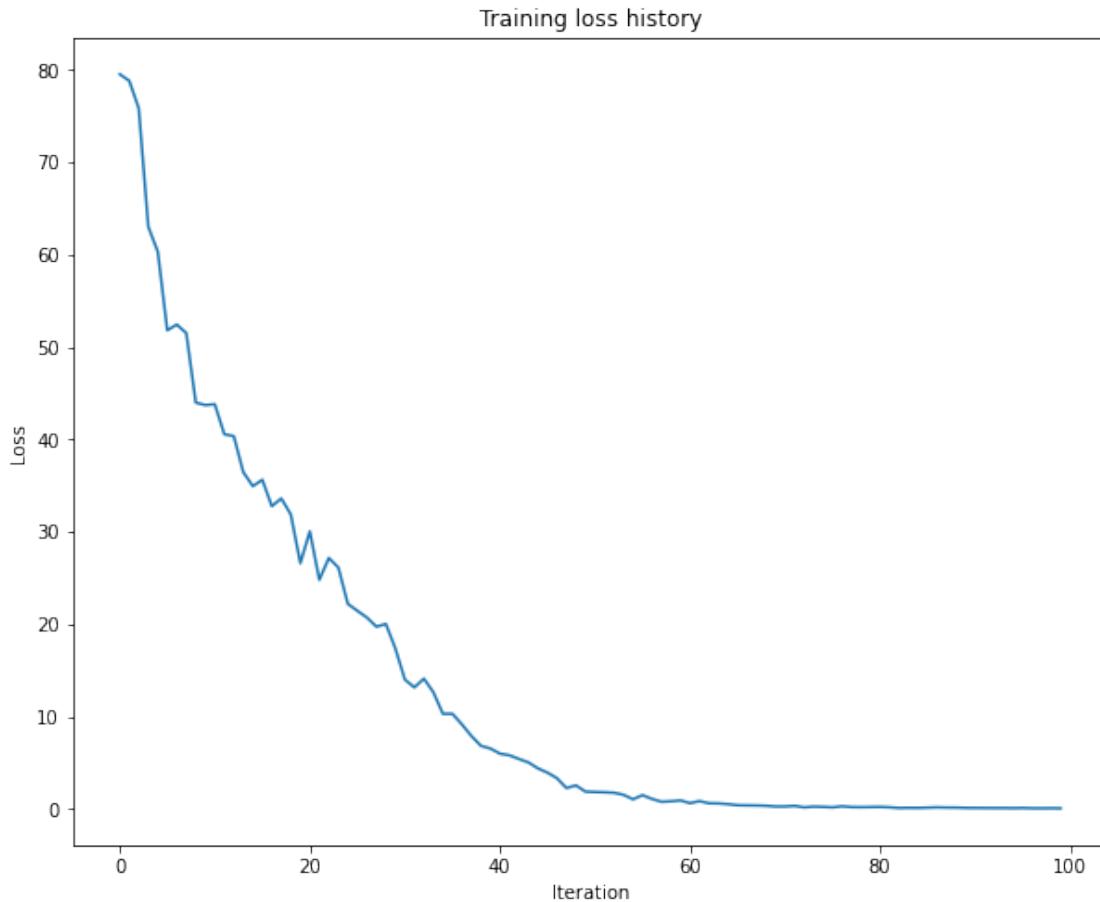
small_lstm_solver = CaptioningSolver(small_lstm_model, small_data,
    update_rule='adam',
    num_epochs=50,
    batch_size=25,
    optim_config={
        'learning_rate': 5e-3,
    },
    lr_decay=0.995,
    verbose=True, print_every=10,
)

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()
```

```
base dir
/Users/luciabouza/FING/DLVis/assignment3_jupyter/cs231n/datasets/coco_captioning
(Iteration 1 / 100) loss: 79.551150
(Iteration 11 / 100) loss: 43.829102
(Iteration 21 / 100) loss: 30.062625
(Iteration 31 / 100) loss: 14.020130
(Iteration 41 / 100) loss: 6.004853
(Iteration 51 / 100) loss: 1.849936
(Iteration 61 / 100) loss: 0.643601
(Iteration 71 / 100) loss: 0.286787
```

```
(Iteration 81 / 100) loss: 0.237477  
(Iteration 91 / 100) loss: 0.126900
```



Print final training loss. You should see a final loss of less than 0.5.

```
[13]: print('Final loss: ', small_lstm_solver.loss_history[-1])
```

```
Final loss:  0.08131789979690769
```

11 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).

```
[14]: 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

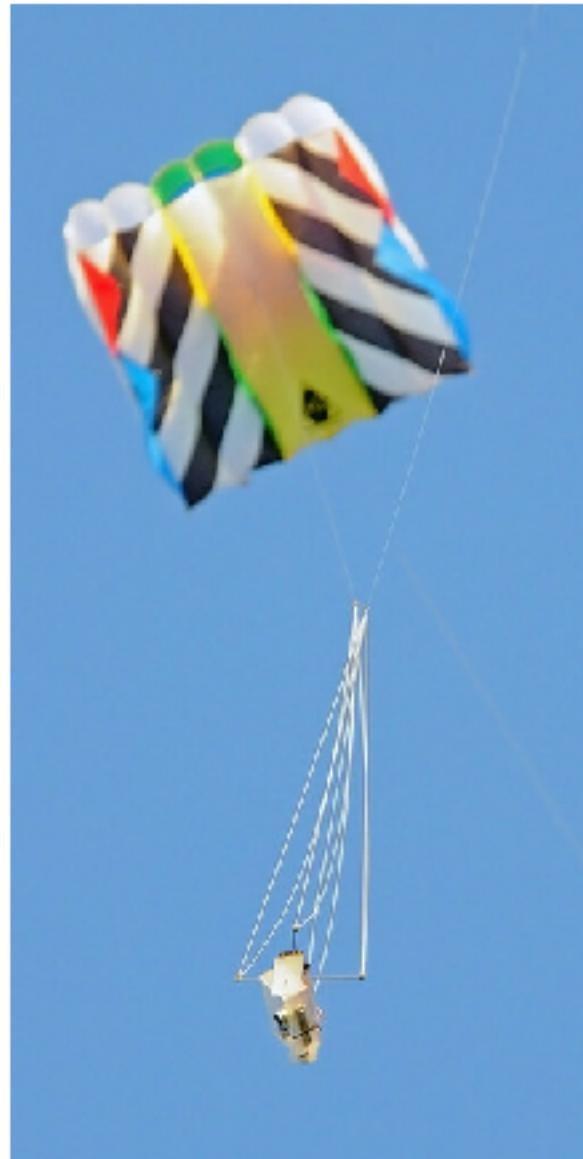
a man standing on the side of a road with bags of luggage <END>
 GT:<START> a man standing on the side of a road with bags of luggage <END>



train

a man <UNK> with a bright colorful kite <END>

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



val

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



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



[]:

Network Visualization-TensorFlow

December 4, 2020

1 Network Visualization (TensorFlow)

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

```
[1]: # As usual, a bit of setup
import sys
sys.path.append('../')
import time, os, json
import numpy as np
import matplotlib.pyplot as plt
import tensorflow as tf

from cs231n.classifiers.squeezenet import SqueezeNet
from cs231n.data_utils import load_tiny_imagenet
from cs231n.image_utils import preprocess_image, deprocess_image
```

```

from cs231n.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
# autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

```

2 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.

We have ported the PyTorch SqueezeNet model to TensorFlow; see: `cs231n/classifiers/squeeze.py` for the model architecture.

To use SqueezeNet, you will need to first **download the weights** by descending into the `cs231n/datasets` directory and running `get_squeeze_tf.sh`. Note that if you ran `get_assignment3_data.sh` then SqueezeNet will already be downloaded.

Once you've downloaded the SqueezeNet model, we can load it into a new TensorFlow session:

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

NOTE: Ignore Tensorflow warnings from the cell below

```

[2]: SAVE_PATH = None

# Local
SAVE_PATH = 'cs231n/datasets/squeeze.ckpt'

# Colab
#SAVE_PATH = '/content/drive/My Drive/{}{}'.format(FOLDERNAME, 'cs231n/
# datasets/squeeze.ckpt')

assert SAVE_PATH is not None, "[!] Choose path to squeeze.ckpt"

if not os.path.exists(SAVE_PATH + ".index"):
    raise ValueError("You need to download SqueezeNet!")

```

```
model = SqueezeNet()
status = model.load_weights(SAVE_PATH)

model.trainable = False
```

WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<tensorflow.python.keras.layers.convolutional.Conv2D object at 0x7fabf43ef520> and <tensorflow.python.keras.layers.pooling.MaxPooling2D object at 0x7fabf43ef340>).

WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf43f92b0> and <cs231n.classifiers.squeezeNet.Fire object at 0x7fabf3e3e640>).
WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf3e3e640> and <tensorflow.python.keras.layers.pooling.MaxPooling2D object at 0x7fabf43f93a0>).
WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf463c490> and <cs231n.classifiers.squeezeNet.Fire object at 0x7fabf502b910>).
WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf502b910> and <tensorflow.python.keras.layers.pooling.MaxPooling2D object at 0x7fabf5031eb0>).
WARNING:tensorflow:Inconsistent references when loading the checkpoint into this

object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf50352b0> and
<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf5039730>).

WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf5039730> and
<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf503ecd0>).

WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf503ecd0> and
<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf50472b0>).

WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<cs231n.classifiers.squeezeNet.Fire object at 0x7fabf50472b0> and
<tensorflow.python.keras.layers.convolutional.Conv2D object at 0x7fabf504abe0>).
WARNING:tensorflow:Inconsistent references when loading the checkpoint into this object graph. Either the Trackable object references in the Python program have changed in an incompatible way, or the checkpoint was generated in an incompatible program.

Two checkpoint references resolved to different objects
(<tensorflow.python.keras.layers.convolutional.Conv2D object at 0x7fabf504abe0>
and <tensorflow.python.keras.layers.pooling.AveragePooling2D object at
0x7fabf504adc0>).

2.1 Load some ImageNet images

We have provided a few example images from the validation set of the ImageNet ILSVRC 2012 Classification dataset. To download these images, descend into `cs231n/datasets/` and run `get_imagenet_val.sh`.

Since they come from the validation set, our pretrained model did not see these images during

training.

Run the following cell to visualize some of these images, along with their ground-truth labels.

```
[3]: from cs231n.data_utils import load_imagenet_val
X_raw, y, class_names = load_imagenet_val(num=5)

plt.figure(figsize=(12, 6))
for i in range(5):
    plt.subplot(1, 5, i + 1)
    plt.imshow(X_raw[i])
    plt.title(class_names[y[i]])
    plt.axis('off')
plt.gcf().tight_layout()
```



2.2 Preprocess images

The input to the pretrained model is expected to be normalized, so we first preprocess the images by subtracting the pixelwise mean and dividing by the pixelwise standard deviation.

```
[4]: X = np.array([preprocess_image(img) for img in X_raw])
```

3 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 $(H, W, 3)$ then this gradient will also have shape $(H, W, 3)$; 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.

Open the file `cs231n/classifiers/squeezezenet.py` and read the code to make sure you understand how to use the model. You will have to use `tf.GradientTape()` to compute gradients with respect to the pixels of the image. In particular, it will be very useful to read this [section](#) for better understanding.

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

Implement `compute_saliency_maps` function inside `cs231n/net_visualization_tensorflow.py`

3.0.1 Hint: Tensorflow `gather_nd` method

Recall in Assignment 1 you needed to select one element from each row of a matrix; if `s` is an numpy array of shape (N, C) and `y` is a numpy array of shape $(N,)$ containing integers $0 \leq y[i] < C$, then `s[np.arange(N), y]` is a numpy array of shape $(N,)$ which selects one element from each element in `s` using the indices in `y`.

In Tensorflow you can perform the same operation using the `gather_nd()` method. If `s` is a Tensor of shape (N, C) and `y` is a Tensor of shape $(N,)$ containing longs in the range $0 \leq y[i] < C$, then

```
tf.gather_nd(s, tf.stack((tf.range(N), y), axis=1))
```

will be a Tensor of shape $(N,)$ containing one entry from each row of `s`, selected according to the indices in `y`.

You can also read the documentation for the [gather_nd method](#).

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

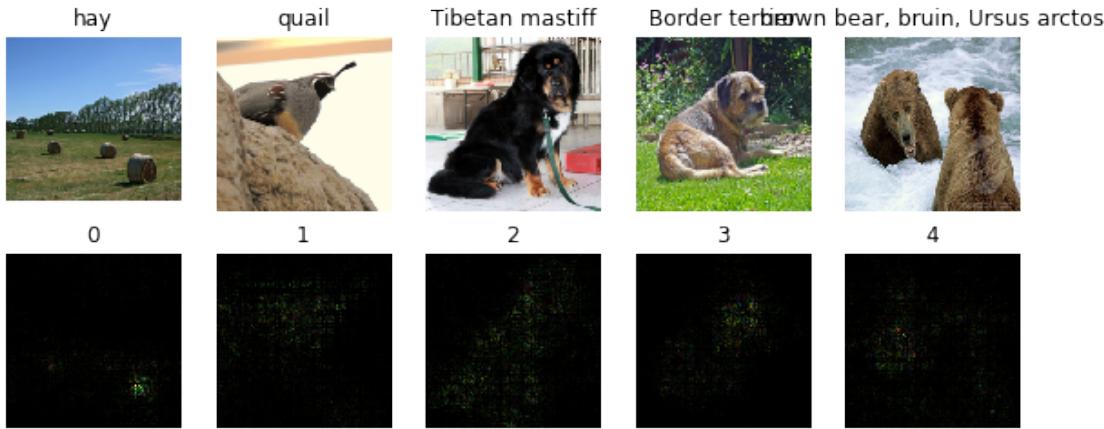
```
[5]: from cs231n.net_visualization_tensorflow import compute_saliency_maps
def show_saliency_maps(X, y, mask):
    mask = np.asarray(mask)
    Xm = X[mask]
    ym = y[mask]

    saliency = compute_saliency_maps(Xm, ym, model)

    for i in range(mask.size):
        plt.subplot(2, mask.size, i + 1)
        plt.imshow(deprocess_image(Xm[i]))
        plt.axis('off')
        plt.title(class_names[ym[i]])
    plt.subplot(2, mask.size, mask.size + i + 1)
    plt.title(mask[i])
    plt.imshow(saliency[i], cmap=plt.cm.hot)
    plt.axis('off')
    plt.gcf().set_size_inches(10, 4)
    plt.show()

mask = np.arange(5)
show_saliency_maps(X, y, mask)
```

Clipping input data to the valid range for `imshow` with RGB data ([0..1] for floats or [0..255] for integers).



4 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:

La respuesta no es correcta. Si se desea maximizar el score de la clase correcta, el método de optimización es ascenso por gradiente, eso si es correcto. Los Saliency maps representan el valor absoluto de los gradientes de los scores correctos con respecto a las imágenes. Si en el paso de actualización usáramos los Saliency maps, en vez de los gradientes, no estaríamos teniendo en cuenta el signo de la dirección de ascenso, por lo que no podemos asegurar acercarnos y converger a un máximo.

5 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

Implement `make_fooling_image` function inside `cs231n/net_visualization_tensforflow.py`

Run the following to generate a fooling image. You should ideally see at first glance no major difference between the original and fooling images, and the network should now make an incorrect prediction on the fooling one. However you should see a bit of random noise if you look at the 10x magnified difference between the original and fooling images. Feel free to change the `idx` variable to explore other images.

```
[6]: from cs231n.net_visualization_tensorflow import make_fooling_image

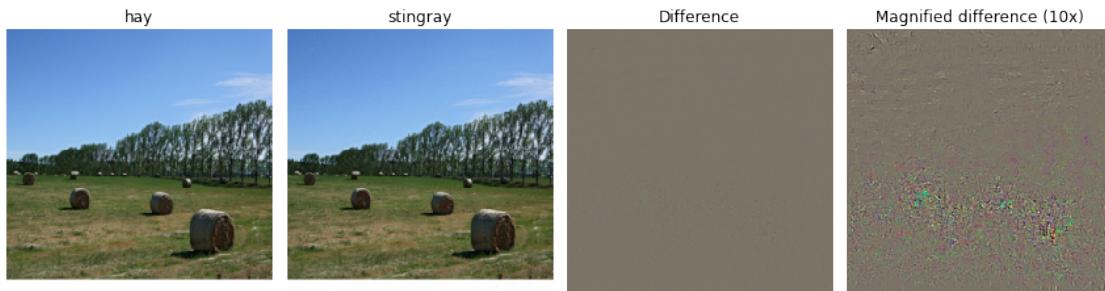
idx = 0
Xi = X[idx][None]
target_y = 6
X_fooling = make_fooling_image(Xi, target_y, model)

# Make sure that X_fooling is classified as y_target
scores = model(X_fooling)
assert tf.math.argmax(scores[0]).numpy() == target_y, 'The network is not\u2192fooled!'

# Show original image, fooling image, and difference
orig_img = deprocess_image(Xi[0])
fool_img = deprocess_image(X_fooling[0])
plt.figure(figsize=(12, 6))

# Rescale
plt.subplot(1, 4, 1)
plt.imshow(orig_img)
plt.axis('off')
plt.title(class_names[y[idx]])
plt.subplot(1, 4, 2)
plt.imshow(fool_img)
plt.title(class_names[target_y])
plt.axis('off')
plt.subplot(1, 4, 3)
plt.title('Difference')
plt.imshow(deprocess_image((Xi-X_fooling)[0]))
plt.axis('off')
plt.subplot(1, 4, 4)
plt.title('Magnified difference (10x)')
plt.imshow(deprocess_image(10 * (Xi-X_fooling)[0]))
plt.axis('off')
plt.gcf().tight_layout()
```

la engaño!



6 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.

[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 `cs231n/net_visualization_tensorflow.py` complete the implementation of the `image_visualization_update_step` used in the `create_class_visualization` function below. Once you have completed that implementation, run the following cells to generate an image of a Tarantula:

```
[7]: from cs231n.net_visualization_tensorflow import*
      ↪class_visualization_update_step, jitter, blur_image

def create_class_visualization(target_y, model, **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
```

```

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 jitter the image as an implicit regularizer
- show_every: How often to show the intermediate result
"""

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)

# We use a single image of random noise as a starting point
X = 255 * np.random.rand(224, 224, 3)
X = preprocess_image(X)[None]

loss = None # scalar loss
grad = None # gradient of loss with respect to model.image, same size as
            # model.image

X = tf.Variable(X)
for t in range(num_iterations):
    # Randomly jitter the image a bit; this gives slightly nicer results
    ox, oy = np.random.randint(0, max_jitter, 2)
    X = jitter(X, ox, oy)
    X = class_visualization_update_step(X, model, target_y, l2_reg,
                                        learning_rate)
    # Undo the jitter
    X = jitter(X, -ox, -oy)
    # As a regularizer, clip and periodically blur

    X = tf.clip_by_value(X, -SQUEEZENET_MEAN/SQUEEZENET_STD, (1.0 -
            SQUEEZENET_MEAN)/SQUEEZENET_STD)
    if t % blur_every == 0:
        X = blur_image(X, sigma=0.5)

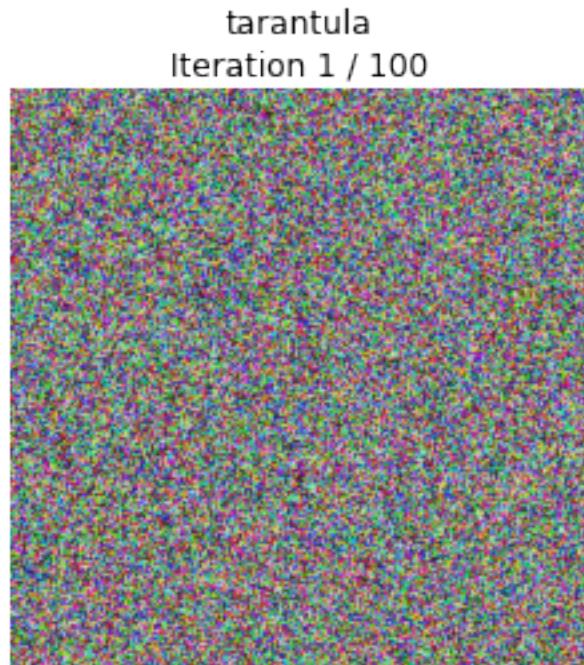
    # Periodically show the image
    if t == 0 or (t + 1) % show_every == 0 or t == num_iterations - 1:
        plt.imshow(deprocess_image(X[0]))
        class_name = class_names[target_y]
        plt.title('%s\nIteration %d / %d' % (class_name, t + 1,
                                              num_iterations))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

```

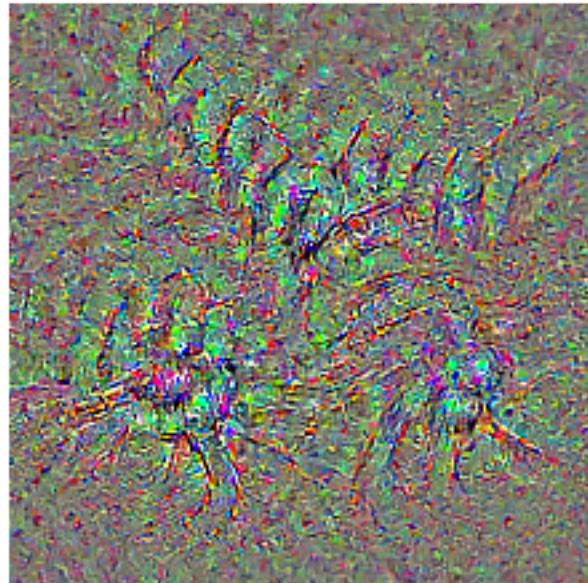
```
    plt.axis('off')
    plt.show()
return X
```

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

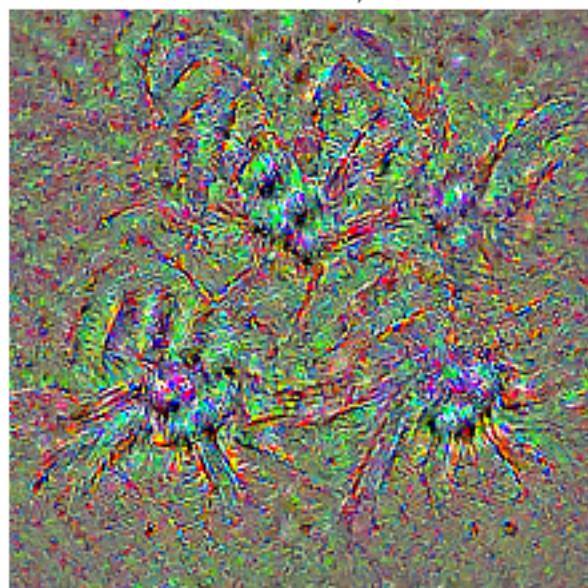
```
[8]: target_y = 76 # Tarantula
out = create_class_visualization(target_y, model)
```



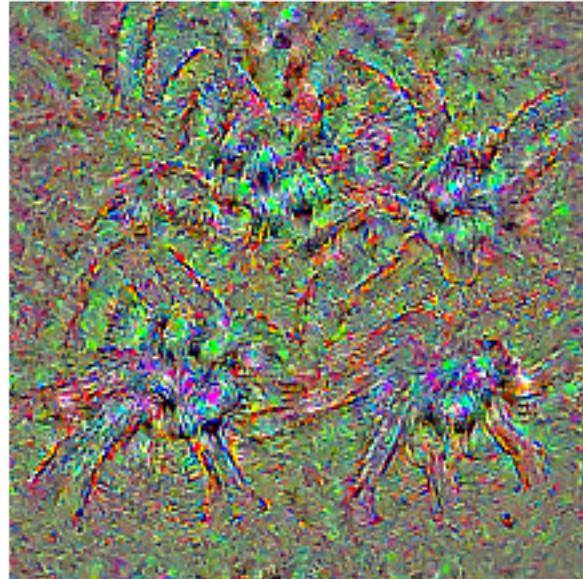
tarantula
Iteration 25 / 100



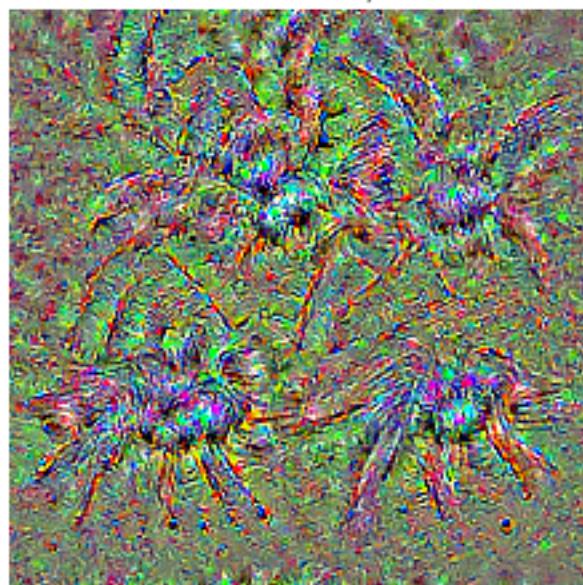
tarantula
Iteration 50 / 100



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.

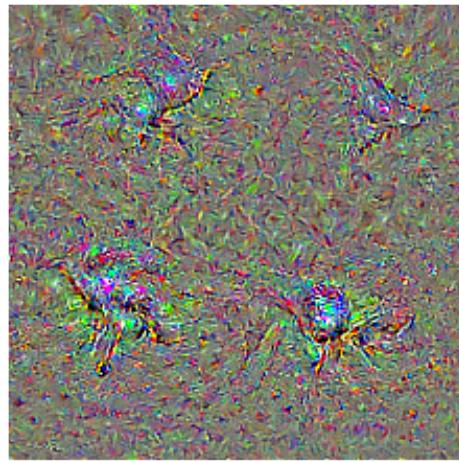
```
[9]: target_y = np.random.randint(1000)
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
print(class_names[target_y])
X = create_class_visualization(target_y, model)
```

platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus

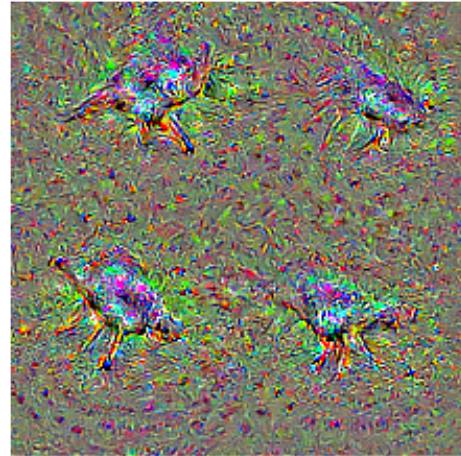
platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus
Iteration 1 / 100



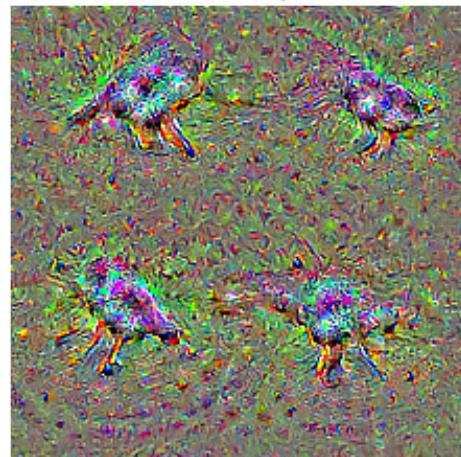
platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus
Iteration 25 / 100



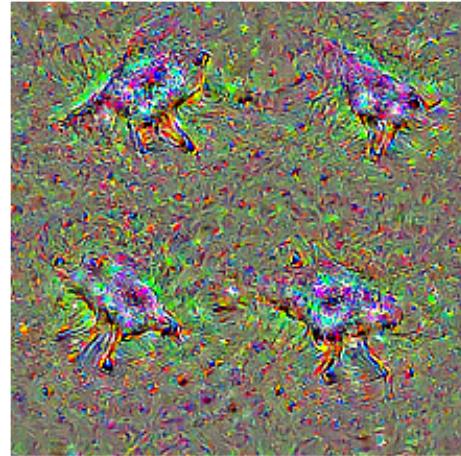
platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus
Iteration 50 / 100



platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus
Iteration 75 / 100



platypus, duckbill, duckbilled platypus, duck-billed platypus, Ornithorhynchus anatinus
Iteration 100 / 100



[]:

StyleTransfer- TensorFlow

December 4, 2020

1 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:



1.1 Part 0: Setup

```
[ ]: import os
import numpy as np
import matplotlib.pyplot as plt
import tensorflow as tf

# Helper functions to deal with image preprocessing
from cs231n.image_utils import load_image, preprocess_image, deprocess_image
from cs231n.classifiers.squeezenet import SqueezeNet

%matplotlib inline
%load_ext autoreload
%autoreload 2
```

```

# Older versions of scipy.misc.imresize yield different results
# from newer versions, so we check to make sure scipy is up to date.
def check_scipy():
    import scipy
    version = scipy.__version__.split('.')
    if int(version[0]) < 1:
        assert int(version[1]) >= 16, "You must install SciPy >= 0.16.0 to"
    ↪complete this notebook."
check_scipy()

```

Load the pretrained SqueezeNet model. This model has been ported from PyTorch, see `cs231n/classifiers/squeezeenet.py` for the model architecture.

To use SqueezeNet, you will need to first **download the weights** by descending into the `cs231n/datasets` directory and running `get_squeezeenet_tf.sh`. Note that if you ran `get_assignment3_data.sh` then SqueezeNet will already be downloaded.

NOTE: Ignore Tensorflow warnings in cells below

```

[ ]: # Load pretrained SqueezeNet model
SAVE_PATH = None

# Local
# SAVE_PATH = 'cs231n/datasets/squeezeenet.ckpt'

# Colab
#SAVE_PATH = '/content/drive/My Drive/{}{}'.format(FOLDERNAME, 'cs231n/'
    ↪datasets/squeezeenet.ckpt')

assert SAVE_PATH is not None, "[!] Choose path to squeezeenet.ckpt"

if not os.path.exists(SAVE_PATH + ".index"):
    raise ValueError("You need to download SqueezeNet!")

CHECKS_PATH = SAVE_PATH.replace('cs231n/datasets/squeezeenet.ckpt', 'style-transfer-checks-tf.npz')
STYLES_FOLDER = CHECKS_PATH.replace('style-transfer-checks-tf.npz', 'styles')

```

```

[ ]: model=SqueezeNet()
model.load_weights(SAVE_PATH)
model.trainable=False

# Load data for testing

```

```

content_img_test = preprocess_image(load_image('%s/tubingen.jpg' %_
    ↪(STYLES_FOLDER), size=192)) [None]
style_img_test = preprocess_image(load_image('%s/starry_night.jpg' %_
    ↪(STYLES_FOLDER), size=192)) [None]
answers = np.load(CHECKS_PATH)

```

2 Part 1: 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.

2.1 Part 1A: Content loss

We can generate an image that reflects the content of one image and the style of another by incorporating both in our loss function. We want to penalize deviations from the content of the content image and deviations from the style of the style image. We can then use this hybrid loss function to perform gradient descent **not on the parameters** of the model, but instead **on the pixel values** of our original image.

Let's first write the content loss function. Content loss measures how much the feature map of the generated image differs from the feature map of the source image. We only care about the content representation of one layer of the network (say, layer ℓ), that has feature maps $A^\ell \in \mathbb{R}^{1 \times H_\ell \times W_\ell \times C_\ell}$. C_ℓ is the number of filters/channels in layer ℓ , H_ℓ and W_ℓ are the height and width. We will work with reshaped versions of these feature maps that combine all spatial positions into one dimension. Let $F^\ell \in \mathbb{R}^{M_\ell \times C_\ell}$ be the feature map for the current image and $P^\ell \in \mathbb{R}^{M_\ell \times C_\ell}$ be the feature map for the content source image where $M_\ell = H_\ell \times W_\ell$ is the number of elements in each feature map. Each row of F^ℓ or P^ℓ represents the vectorized activations of a particular filter, convolved over all positions of the image. Finally, let w_c be the weight of the content loss term in the loss function.

Then the content loss is given by:

$$L_c = w_c \times \sum_{i,j} (F_{ij}^\ell - P_{ij}^\ell)^2$$

Implement `content_loss` in `cs231n/style_transfer_tensorflow.py`

Test your content loss. The error should be less than 1e-8.

```
[ ]: from cs231n.style_transfer_tensorflow import content_loss, extract_features,_
    ↪rel_error
def content_loss_test(correct):
    content_layer = 2
    content_weight = 6e-2
    c_feats = extract_features(content_img_test, model)[content_layer]
    bad_img = tf.zeros(content_img_test.shape)
    feats = extract_features(bad_img, model)[content_layer]
    student_output = content_loss(content_weight, c_feats, feats)
    error = rel_error(correct, student_output)
    print('Maximum error is {:.3f}'.format(error))
```

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

2.2 Part 1B: Style loss

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

First, compute the Gram matrix G which represents the correlations between the values in each channel of the feature map (i.e. the “responses” of the filter responsible for that channel), where F is as above. The Gram matrix is an approximation of the covariance matrix – it tells us how every channel’s values (i.e. that filter’s activations) correlate with every other channel’s values. If we have C channels, matrix G will be of shape (C, C) to capture these correlations.

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^ℓ of shape (C_ℓ, H_ℓ, W_ℓ) , we can flatten the height and width dimensions so they’re just 1 dimension $M_\ell = H_\ell \times W_\ell$: the new shape of F^ℓ is (C_ℓ, M_ℓ) . Then, the Gram matrix has shape (C_ℓ, C_ℓ) where each element is given by the equation:

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

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

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

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

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

Begin by implementing the Gram matrix computation function `gram_matrix` inside `cs231n\style_transfer_tensorflow.py`:

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

```
[ ]: from cs231n.style_transfer_tensorflow import gram_matrix

def gram_matrix_test(correct):
    gram = gram_matrix(extract_features(style_img_test, model)[4]) ### 4 ↳ instead of 5 - second MaxPooling layer
    error = rel_error(correct, gram)
```

```

    print('Maximum error is {:.3f}'.format(error))

gram_matrix_test(answers['gm_out'])

```

Next, put it together and implement the style loss function `style_loss` in `cs231n/style_transfer_tensorflow.py`

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

```

[ ]: from cs231n.style_transfer_tensorflow import style_loss

def style_loss_test(correct):
    style_layers = [0, 3, 5, 6]
    style_weights = [300000, 1000, 15, 3]

    c_feats = extract_features(content_img_test, model)
    feats = extract_features(style_img_test, model)
    style_targets = []
    for idx in style_layers:
        style_targets.append(gram_matrix(feats[idx]))

    s_loss = style_loss(c_feats, style_layers, style_targets, style_weights)
    error = rel_error(correct, s_loss)
    print('Error is {:.3f}'.format(error))

style_loss_test(answers['sl_out'])

```

2.3 Part 1C: 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 regularization 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 `cs231/style_transfer_tensorflow.py`, fill in the definition for the TV loss term in `tv_loss`. To receive full credit, your implementation should not have any loops.

Test your TV loss implementation. Error should be less than 0.001.

```

[ ]: from cs231n.style_transfer_tensorflow import tv_loss
from inspect import getsourcelines
import re
def tv_loss_test(correct):
    tv_weight = 2e-2
    t_loss = tv_loss(content_img_test, tv_weight)

```

```

error = rel_error(correct, t_loss)
print('Error is {:.4f}'.format(error))

lines, _ = getsourcelines(tv_loss)
used_loop = any(bool(re.search(r"for \S* in", line)) for line in lines)
if used_loop:
    print("WARNING!!!! - Your implementation of tv_loss contains a loop! To\u202a
→receive full credit, your implementation should not have any loops")

tv_loss_test(answers['tv_out'])

```

3 Part 2: Style Transfer

Lets put it all together and make some beautiful images! The `style_transfer` function below combines all the losses you coded up above and optimizes for an image that minimizes the total loss.

```
[ ]: def style_transfer(content_image, style_image, image_size, style_size,✉
→content_layer, content_weight,
                      style_layers, style_weights, tv_weight, init_random = False):
    """Run style transfer!

    Inputs:
    - content_image: filename of content image
    - style_image: filename of style image
    - image_size: size of smallest image dimension (used for content loss and✉
→generated image)
    - style_size: size of smallest style image dimension
    - content_layer: layer to use for content loss
    - content_weight: weighting on content loss
    - style_layers: list of layers to use for style loss
    - style_weights: list of weights to use for each layer in style_layers
    - tv_weight: weight of total variation regularization term
    - init_random: initialize the starting image to uniform random noise
    """

    # Extract features from the content image
    content_img = preprocess_image(load_image(content_image, size=image_size))
    feats = extract_features(content_img[None], model)
    content_target = feats[content_layer]

    # Extract features from the style image
    style_img = preprocess_image(load_image(style_image, size=style_size))
    s_feats = extract_features(style_img[None], model)
    style_targets = []
    # Compute list of TensorFlow Gram matrices
    for idx in style_layers:
```

```

    style_targets.append(gram_matrix(s_feats[idx]))

# Set up optimization hyperparameters
initial_lr = 3.0
decayed_lr = 0.1
decay_lr_at = 180
max_iter = 200

step = tf.Variable(0, trainable=False)
boundaries = [decay_lr_at]
values = [initial_lr, decayed_lr]
learning_rate_fn = tf.keras.optimizers.schedules.
    ↪PiecewiseConstantDecay(boundaries, values)

# Later, whenever we perform an optimization step, we pass in the step.
learning_rate = learning_rate_fn(step)

optimizer = tf.keras.optimizers.Adam(learning_rate=learning_rate)

# Initialize the generated image and optimization variables

f, axarr = plt.subplots(1,2)
axarr[0].axis('off')
axarr[1].axis('off')
axarr[0].set_title('Content Source Img.')
axarr[1].set_title('Style Source Img.')
axarr[0].imshow(deprocess_image(content_img))
axarr[1].imshow(deprocess_image(style_img))
plt.show()
plt.figure()

# Initialize generated image to content image
if init_random:
    initializer = tf.random_uniform_initializer(0, 1)
    img = initializer(shape=content_img[None].shape)
    img_var = tf.Variable(img)
    print("Intializing randomly.")
else:
    img_var = tf.Variable(content_img[None])
    print("Initialzing with content image.")

for t in range(max_iter):
    with tf.GradientTape() as tape:
        tape.watch(img_var)
        feats = extract_features(img_var, model)
        # Compute loss

```

```

        c_loss = content_loss(content_weight, feats[content_layer], ↴
    ↪content_target)
        s_loss = style_loss(feats, style_layers, style_targets, ↴
    ↪style_weights)
        t_loss = tv_loss(img_var, tv_weight)
        loss = c_loss + s_loss + t_loss
    # Compute gradient
    grad = tape.gradient(loss, img_var)
    optimizer.apply_gradients([(grad, img_var)])

    img_var.assign(tf.clip_by_value(img_var, -1.5, 1.5))

    if t % 100 == 0:
        print('Iteration {}'.format(t))
        plt.imshow(deprocess_image(img_var[0].numpy(), rescale=True))
        plt.axis('off')
        plt.show()
    print('Iteration {}'.format(t))
    plt.imshow(deprocess_image(img_var[0].numpy(), rescale=True))
    plt.axis('off')
    plt.show()

```

3.1 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.

```
[ ]: # Composition VII + Tubingen
params1 = {
    'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image' : '%s/composition_vii.jpg' % (STYLES_FOLDER),
    'image_size' : 192,
    'style_size' : 512,
    'content_layer' : 2,
    'content_weight' : 5e-2,
    'style_layers' : (0, 3, 5, 6),
    'style_weights' : (20000, 500, 12, 1),
    'tv_weight' : 5e-2
}

style_transfer(**params1)
```

```
[ ]: # Scream + Tubingen
params2 = {
    'content_image': '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image': '%s/the_scream.jpg' % (STYLES_FOLDER),
    'image_size': 192,
    'style_size': 224,
    'content_layer': 2,
    'content_weight': 3e-2,
    'style_layers': [0, 3, 5, 6],
    'style_weights': [200000, 800, 12, 1],
    'tv_weight': 2e-2
}

style_transfer(**params2)
```

```
[ ]: # Starry Night + Tubingen
params3 = {
    'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image' : '%s/starry_night.jpg' % (STYLES_FOLDER),
    'image_size' : 192,
    'style_size' : 192,
    'content_layer' : 2,
    'content_weight' : 6e-2,
    'style_layers' : [0, 3, 5, 6],
    'style_weights' : [300000, 1000, 15, 3],
    'tv_weight' : 2e-2
}

style_transfer(**params3)
```

4 Part 3: 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 “[Understanding Deep Image Representations by Inverting Them](#)” 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

```
[ ]: # Feature Inversion -- Starry Night + Tubingen
params_inv = {
    'content_image' : '%s/tubingen.jpg' % (STYLES_FOLDER),
    'style_image' : '%s/starry_night.jpg' % (STYLES_FOLDER),
    'image_size' : 192,
    'style_size' : 192,
    'content_layer' : 2,
    'content_weight' : 6e-2,
    'style_layers' : [0, 3, 5, 6],
    '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)
```

```
[ ]:
```

Generative_Adversarial_Networks_TF

December 4, 2020

1 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 repertoire, 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.

1.0.1 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:

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

where $x \sim p_{\text{data}}$ are samples from the input data, $z \sim p(z)$ 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 alleviate 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:

$$\underset{G}{\text{maximize}} \mathbb{E}_{z \sim p(z)} [\log D(G(z))]$$

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

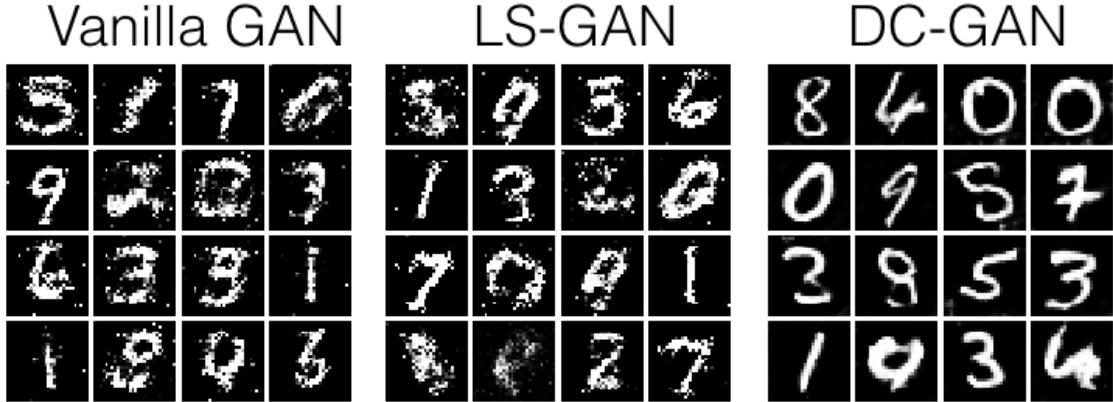
$$\underset{D}{\text{maximize}} \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D(G(z)))]$$

1.0.2 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 stability 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](#)). Variational autoencoders combine neural networks with variational 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.

Example pictures of what you should expect (yours might look slightly different):



1.1 Setup

```
[ ]: import tensorflow as tf
import numpy as np
import os

import matplotlib.pyplot as plt
import matplotlib.gridspec as gridspec

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
# autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

# A bunch of utility functions

def show_images(images):
    images = np.reshape(images, [images.shape[0], -1]) # images reshape to (batch_size, D)
    sqrt_n = int(np.ceil(np.sqrt(images.shape[0])))
    sqrt_m = int(np.ceil(np.sqrt(images.shape[1])))

    fig = plt.figure(figsize=(sqrt_n, sqrt_n))
    gs = gridspec.GridSpec(sqrt_n, sqrt_n)
    gs.update(wspace=0.05, hspace=0.05)

    for i, img in enumerate(images):
        ax = plt.subplot(gs[i])
        plt.axis('off')
        ax.set_xticklabels([])
        ax.set_yticklabels([])
        ax.set_aspect('equal')
        plt.imshow(img.reshape([sqrt_m,sqrt_m]))
    return

from cs231n.gan_tf import preprocess_img, deprocess_img, rel_error,
count_params, MNIST

NOISE_DIM = 96
```

```
[ ]: # Colab users only  
%cd drive/My\ Drive/$FOLDERNAME/  
%cp -r gan-checks-tf.npz /content/  
%cd /content/  
  
[ ]: answers = np.load('gan-checks-tf.npz')
```

1.2 Dataset

GANs are notoriously finicky with hyperparameters, and also require many training epochs. In order to make this assignment approachable without a GPU, we will be working on the MNIST dataset, which is 60,000 training and 10,000 test images. Each picture contains a centered image of white digit on black background (0 through 9). This was one of the first datasets used to train convolutional neural networks and it is fairly easy – a standard CNN model can easily exceed 99% accuracy.

Heads-up: Our MNIST wrapper returns images as vectors. That is, they're size (batch, 784). If you want to treat them as images, we have to resize them to (batch,28,28) or (batch,28,28,1). They are also type np.float32 and bounded [0,1].

```
[ ]: # show a batch  
mnist = MNIST(batch_size=16)  
show_images(mnist.X[:16])
```

1.3 LeakyReLU

In the cell below, you should implement a LeakyReLU. See the [class notes](#) (where alpha is small number) or equation (3) in [this paper](#). LeakyReLUs keep ReLU units from dying and are often used in GAN methods (as are maxout units, however those increase model size and therefore are not used in this notebook).

HINT: You should be able to use `tf.maximum`

Implement `leaky_relu` in `cs231n/gan_tf.py`

Test your leaky ReLU implementation. You should get errors < 1e-10

```
[ ]: from cs231n.gan_tf import leaky_relu  
  
def test_leaky_relu(x, y_true):  
    y = leaky_relu(tf.constant(x))  
    print('Maximum error: %g' % rel_error(y_true, y))  
  
test_leaky_relu(answers['lrelu_x'], answers['lrelu_y'])
```

1.4 Random Noise

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

Implement `sample_noise` in `cs231n/gan_tf.py`

Make sure noise is the correct shape and type:

```
[ ]: from cs231n.gan_tf import sample_noise

def test_sample_noise():
    batch_size = 3
    dim = 4
    z = sample_noise(batch_size, dim)
    # Check z has the correct shape
    assert z.get_shape().as_list() == [batch_size, dim]
    # Make sure z is a Tensor and not a numpy array
    assert isinstance(z, tf.Tensor)
    # Check that we get different noise for different evaluations
    z1 = sample_noise(batch_size, dim)
    z2 = sample_noise(batch_size, dim)
    assert not np.array_equal(z1, z2)
    # Check that we get the correct range
    assert np.all(z1 >= -1.0) and np.all(z1 <= 1.0)
    print("All tests passed!")

test_sample_noise()
```

1.5 Discriminator

Our first step is to build a discriminator. **Hint:** You should use the layers in `tf.keras.layers` to build the model. All fully connected layers should include bias terms. For initialization, just use the default initializer used by the `tf.keras.layers` functions.

Architecture: * Fully connected layer with input size 784 and output size 256 * LeakyReLU with alpha 0.01 * Fully connected layer with output size 256 * LeakyReLU with alpha 0.01 * Fully connected layer with output size 1

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

Implement `discriminator()` in `cs231n/gan_tf.py`

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

```
[ ]: from cs231n.gan_tf import discriminator

def test_discriminator(true_count=267009, discriminator=discriminator):
    model = discriminator()
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in discriminator. {} instead of {}'.format(cur_count, true_count))
    else:
        print('Correct number of parameters in discriminator.')
```

```
test_discriminator()
```

1.6 Generator

Now to build a generator. You should use the layers in `tf.keras.layers` to construct the model. All fully connected layers should include bias terms. Note that you can use the `tf.nn` module to access activation functions. Once again, use the default initializers for parameters.

Architecture:

- * Fully connected layer with input size `tf.shape(z)[1]` (the number of noise dimensions) and output size 1024
- * ReLU
- * Fully connected layer with output size 1024
- * ReLU
- * Fully connected layer with output size 784
- * TanH (To restrict every element of the output to be in the range [-1,1])

Implement `generator()` in `cs231n/gan_tf.py`

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

```
[ ]: from cs231n.gan_tf import generator

def test_generator(true_count=1858320, generator=generator):
    model = generator(4)
    cur_count = count_params(model)
    if cur_count != true_count:
        print('Incorrect number of parameters in generator. {0} instead of {1}. □')
        ↪Check your architecture.'.format(cur_count, true_count))
    else:
        print('Correct number of parameters in generator.')

test_generator()
```

2 GAN Loss

Compute the generator and discriminator loss. The generator loss is:

$$\ell_G = -\mathbb{E}_{z \sim p(z)} [\log D(G(z))]$$

and the discriminator loss is:

$$\ell_D = -\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] - \mathbb{E}_{z \sim p(z)} [\log (1 - D(G(z)))]$$

Note that these are negated from the equations presented earlier as we will be *minimizing* these losses.

HINTS: Use `tf.ones` and `tf.zeros` to generate labels for your discriminator. Use `tf.keras.losses.BinaryCrossentropy` to help compute your loss function.

Implement `discriminator_loss`, `generator_loss` in `cs231n/gan_tf.py`

Test your GAN loss. Make sure both the generator and discriminator loss are correct. You should see errors less than 1e-8.

```
[ ]: from cs231n.gan_tf import discriminator_loss

def test_discriminator_loss(logits_real, logits_fake, d_loss_true):
    d_loss = discriminator_loss(tf.constant(logits_real),
                                tf.constant(logits_fake))
    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'])
```

```
[ ]: from cs231n.gan_tf import generator_loss

def test_generator_loss(logits_fake, g_loss_true):
    g_loss = generator_loss(tf.constant(logits_fake))
    print("Maximum error in g_loss: %g"%rel_error(g_loss_true, g_loss))

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

3 Optimizing our loss

Make an Adam optimizer with a 1e-3 learning rate, beta1=0.5 to minimize G_loss and D_loss separately. The trick of decreasing beta was shown to be effective in helping GANs converge in the [Improved Techniques for Training GANs](#) paper. In fact, with our current hyperparameters, if you set beta1 to the Tensorflow default of 0.9, there's a good chance your discriminator loss will go to zero and the generator will fail to learn entirely. In fact, this is a common failure mode in GANs; if your D(x) learns too fast (e.g. loss goes near zero), your G(z) is never able to learn. Often D(x) is trained with SGD with Momentum or RMSProp instead of Adam, but here we'll use Adam for both D(x) and G(z).

Implement get_solvers in cs231n/gan_tf.py

```
[ ]: from cs231n.gan_tf import get_solvers
```

4 Training a GAN!

Well that wasn't so hard, was it? After the first epoch, you should see fuzzy outlines, clear shapes as you approach epoch 3, and decent shapes, about half of which will be sharp and clearly recognizable as we pass epoch 5. In our case, we'll simply train D(x) and G(z) with one batch each every iteration. However, papers often experiment with different schedules of training D(x) and G(z), sometimes doing one for more steps than the other, or even training each one until the loss gets "good enough" and then switching to training the other.

If you are a Colab user, it is recommended to change colab runtime to GPU.

Train your GAN! This should take about 10 minutes on a CPU, or about 2 minutes on GPU.

```
[ ]: from cs231n.gan_tf import run_a_gan

# Make the discriminator
D = discriminator()

# Make the generator
G = generator()

# Use the function you wrote earlier to get optimizers for the Discriminator
# and the Generator
D_solver, G_solver = get_solvers()

# Run it!
images, final = run_a_gan(D, G, D_solver, G_solver, discriminator_loss,
                           generator_loss)
```

```
[ ]: numIter = 0
for img in images:
    print("Iter: {}".format(numIter))
    show_images(img)
    plt.show()
    numIter += 20
    print()
```

```
[ ]: print('Vanilla GAN Final images')
show_images(final)
plt.show()
```

5 Least Squares GAN

We'll now look at [Least Squares GAN](#), a newer, more stable alternative 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:

$$\ell_G = \frac{1}{2} \mathbb{E}_{z \sim p(z)} [(D(G(z)) - 1)^2]$$

and the discriminator loss:

$$\ell_D = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} [(D(x) - 1)^2] + \frac{1}{2} \mathbb{E}_{z \sim p(z)} [(D(G(z)))^2]$$

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 (`score_real` and `score_fake`).

Implement `ls_discriminator_loss`, `ls_generator_loss` in `cs231n/gan_tf.py`

Test your LSGAN loss. You should see errors less than 1e-8.

```
[ ]: from cs231n.gan_tf import ls_discriminator_loss, ls_generator_loss

def test_lsgan_loss(score_real, score_fake, d_loss_true, g_loss_true):

    d_loss = ls_discriminator_loss(tf.constant(score_real), tf.
→constant(score_fake))
    g_loss = ls_generator_loss(tf.constant(score_fake))
    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'])
```

Create new training steps so we instead minimize the LSGAN loss:

```
[ ]: # Make the discriminator
D = discriminator()

# Make the generator
G = generator()

# Use the function you wrote earlier to get optimizers for the Discriminator
→and the Generator
D_solver, G_solver = get_solvers()

# Run it!
images, final = run_a_gan(D, G, D_solver, G_solver, ls_discriminator_loss,
→ls_generator_loss)
```

```
[ ]: numIter = 0
for img in images:
    print("Iter: {}".format(numIter))
    show_images(img)
    plt.show()
    numIter += 20
    print()
```

```
[ ]: print('LSGAN Final images')
show_images(final)
plt.show()
```

6 Deep 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 as our discriminators and generators.

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. *Be sure to check the dimensions of x and reshape when needed*, fully connected blocks expect [N,D] Tensors while conv2d blocks expect [N,H,W,C] Tensors. Please use `tf.keras.layers` to define the following architecture:

Architecture: * Conv2D: 32 Filters, 5x5, Stride 1, padding 0 * Leaky ReLU(alpha=0.01) * Max Pool 2x2, Stride 2 * Conv2D: 64 Filters, 5x5, Stride 1, padding 0 * 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

Once again, please use biases for all convolutional and fully connected layers, and use the default parameter initializers. Note that a padding of 0 can be accomplished with the ‘VALID’ padding option.

Implement `dc_discriminator` in `cs231n/gan_tf.py`

```
[ ]: from cs231n.gan_tf import dc_discriminator

# model = dc_discriminator()
test_discriminator(1102721, dc_discriminator)
```

Generator For the generator, we will copy the architecture exactly from the [InfoGAN paper](#). See Appendix C.1 MNIST. Please use `tf.keras.layers` for your implementation. You might find the documentation for `tf.keras.layers.Conv2DTranspose` useful. The architecture is as follows.

Architecture: * Fully connected with output size 1024 * ReLU * BatchNorm * Fully connected with output size 7 x 7 x 128 * ReLU * BatchNorm * Resize into Image Tensor of size 7, 7, 128 * Conv2D \hat{T} (transpose): 64 filters of 4x4, stride 2 * ReLU * BatchNorm * Conv2d \hat{T} (transpose): 1 filter of 4x4, stride 2 * TanH

Once again, use biases for the fully connected and transpose convolutional layers. Please use the default initializers for your parameters. For padding, choose the ‘same’ option for transpose convolutions. For Batch Normalization, assume we are always in ‘training’ mode.

Implement `dc_generator` in `cs231n/gan_tf.py`

```
[ ]: from cs231n.gan_tf import dc_generator

test_generator(6595521, generator=dc_generator)
```

We have to recreate our network since we’ve changed our functions.

6.0.1 Train and evaluate a DCGAN

This is the one part of A3 that significantly benefits from using a GPU. It takes 3 minutes on a GPU for the requested five epochs. Or about 50 minutes on a dual core laptop on CPU (feel free

to use 3 epochs if you do it on CPU).

```
[ ]: # Make the discriminator
D = dc_discriminator()

# Make the generator
G = dc_generator()

# Use the function you wrote earlier to get optimizers for the Discriminator and the Generator
D_solver, G_solver = get_solvers()

# Run it!
images, final = run_a_gan(D, G, D_solver, G_solver, discriminator_loss, generator_loss, num_epochs=5)
```

```
[ ]: numIter = 0
for img in images:
    print("Iter: {}".format(numIter))
    show_images(img)
    plt.show()
    numIter += 20
    print()
```

```
[ ]: print('DCGAN Final images')
show_images(final)
plt.show()
```

```
[ ]:
```

6.1 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.

6.1.1 Your answer:

| y_0 | y_1 | y_2 | y_3 | y_4 | y_5 | y_6 |
|-------|-------|-------|-------|-------|-------|-------|
| 1 | | | | | | |
| x_0 | x_1 | x_2 | x_3 | x_4 | x_5 | x_6 |
| 1 | | | | | | |

6.2 INLINE QUESTION 2

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

6.2.1 Your answer:

6.3 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.

6.3.1 Your answer:

[]: