

CS7643: Deep Learning
Fall 2020
HW4 Solutions

Tianxue Hu

October 21, 2020

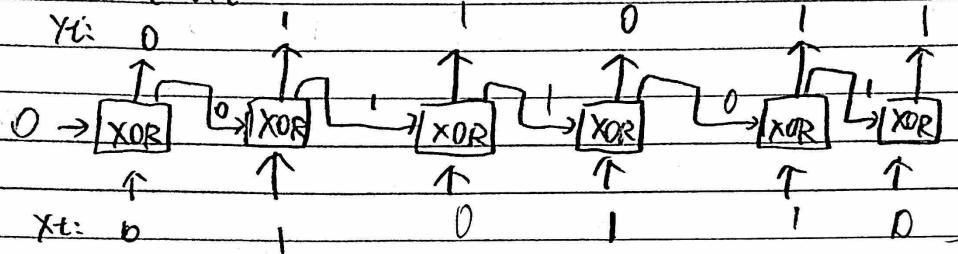
Problem 1

1. Design a sequence parity function with vanilla RNN. we can use XOR
 input $x = [0, 1, 0, 1, 1, 0]$ and output is $[0, 1, 1, 0, 1, 1]$
 the hidden units can be written as:

$$h_0 = 0$$

$$h_t = h_{t-1} \text{ XOR } x_t$$

$$y_t = h_t$$



XOR can be represented as $\text{XOR}(a, b) = \bar{a} \cdot \bar{b} + \bar{a} \cdot b$

$$\text{XOR}(0, 0) = 0$$

$$\text{XOR}(1, 0) = 1$$

$$\text{XOR}(0, 1) = 1$$

$$\text{XOR}(1, 1) = 0$$

In our RNN, we have

y_{t-1}	x_t	y_t
0	0	0
0	1	1
1	1	0
1	0	1

Problem 2

2. LSTM parity function

input x_t, h_{t-1}

Output y_t, \tilde{y}_t

Since in RNN, we apply XOR

$$y_t = (\bar{x}_t \wedge y_{t-1}) \vee (x_t \wedge \bar{y}_{t-1})$$

$$= \bar{x}_t * y_{t-1} + x_t * \bar{y}_{t-1}$$

In LSTM, we have

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

Comparing the two equations, we have

$$c_t = y_t = h_t$$

$$f_t = \bar{x}_t$$

$$i_t = x_t$$

- So we set $c_t = h_t$. to find w_f, b_f .

- We solve.

$$\hat{f}_t = \sigma(w_o \cdot [h_{t-1} \cdot x_t] + b_o) = \bar{x}_t$$

$$\begin{cases} w_f = \begin{bmatrix} 0 \\ 1 \end{bmatrix} & \text{when } x_t=1 \quad \hat{f}_t = \sigma(0) = 0 = \bar{x}_t \\ b_f = 1 & \text{when } x_t=0 \quad \hat{f}_t = \sigma(1) = 1 = \bar{x}_t \end{cases}$$

- Solving w_i, b_i we solve

$$i_t = \sigma(w_i \cdot [h_{t-1} \cdot x_t] + b_i) = x_t$$

$$w_i = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$b_i = 0$$

$$\text{since } x_t=1 \quad \sigma(1) = 1 = x_t$$

$$x_t=0 \quad \sigma(0) = 0 = x_t$$

z cont.

-Solving w.c. b.c. we solve

$$\bar{c}_t = \tanh(w_c [ht_{-1}, x_t] + b_c) = ht_{-1}$$

$$\boxed{w_c = \begin{bmatrix} 1 \\ 0 \end{bmatrix}} \quad b_c = 1 \quad \text{since}$$

$$\bar{c}_t = \tanh(1 - ht_{-1})$$

$$= 1 - ht_{-1} = ht_{-1} \quad \text{where } 1 - ht_{-1} \in [0, 1]$$

-Solving ~~for~~ w₀, b₀, we solve

$$ht = ot * \tanh(\bar{c}_t)$$

$$ot = 1 \quad \text{from } ht = ct \cdot \sin(ht, \dots)$$

$$ot = b(w_0 [ht_{-1}, x_t] + b_0) = 1$$

$$\boxed{w_0 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}}$$

$$\boxed{b_0 = 1}$$

Problem 3

3. When to stop beam search.

The highest scoring beam B_t scores worse than or equal to best_t .

Define $S(B_t)$ be to score of a beam B_t .
Then

$$S(B_t) \leq \text{best}_t$$

Assuming at another time t' that $t' > t$ and

$$S(B_{t'}) \leq S(B_t)$$

as probability value are always ≤ 1
hence

$$S(B_{t'}) \leq (\text{best}_{t'})$$

Proved.

Problem 4

? 4. Exploding Gradient.

$$h_t = w^T h_{t-1}$$

$$h_{t-1} = w^T h_{t-2}$$

$$\vdots$$

$$h_1 = w^T h_0$$

$$\text{Then } h_t = (w^T)^t h_0$$

w could be decomposed as

$$w = P D P^{-1}$$

where D is diagonal matrix of eigenvalues.
 P is the eigen vector of w

Then we can write

$$h_t = ((P D P^{-1})^T)^t h_0$$

$$= ((P^{-1})^T D^T P^T)^t h_0$$

$$= ((P^{-1})^T D P^T)^t h_0$$

$$= (P^{-1})^T D P^T (P^{-1})^T D P^T \dots (P^{-1})^T D P^T h_0$$

$$\text{as } (P^{-1})^T \cdot P^T = I$$

$$= (P^{-1})^T D^t P^T h_0$$

$$\frac{h_t}{h_0} = (P^{-1})^T D^t P^T$$

when $t \geq 0$, if a eigenvalues in $D < 1$, the value in $D^t \rightarrow 0$
 thus vanishing the gradient

if a eigenvalue in $D > 1$, the value in $D^t \rightarrow \infty$, thus
 exploding the gradient

Problem 5

5

$$(a) h_i^{t+1} = g(h_i^t, \text{Agg}(H_{i:t}'))$$

$$= g(h_i^t, \bar{z}_{fj_i}(h_j^t)) / \sum_{j \neq i}$$

where $v_j \in N(v_i)$

$$= g(h_i^t, \sum_{j \in N(v_i)} (v_j^t h_j^t))$$

$$(b) f(h_1^t) = [-2, -2]$$

$$f(h_2^t) = [0, -2]$$

$$f(h_3^t) = [2, 0]$$

$$\text{Agg}(H_{i:t}') = [0.6, 0.2, 0.2] \begin{bmatrix} [-2, 2] \\ [0, -2] \\ [2, 0] \end{bmatrix}$$

$$= [-0.8, 0.8]$$

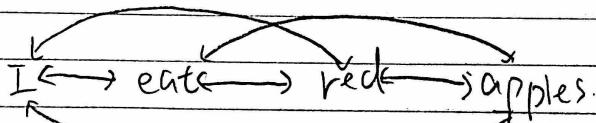
$$h_i^{t+1} = w(h_i^t)^T + \max\{\text{Agg}(H_{i:t}'), 0\}$$

$$= [1, 1] \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \max\{\text{Agg}(H_{i:t}'), 0\}$$

$$= 0 + [-0, 0.8]$$

$$= [0, 0.8]$$

(c)



G has 12 edges

5(d)

$$(17) h_i^{l+1} = \text{Attention}(Q^l h_i^l, K^l h_i^l, V^l h_i^l)$$

$$= \sum_{j \in N(i)} (\text{softmax}_j(Q^l h_i^l \cdot K^l h_j^l) V^l h_j^l)$$

(15) can be

$$h_i^{l+1} = g(h_i^l), \quad \sum_{j \in N(i)} (V^l h_j^l)$$

in this case, we can write w_{ij} as

$$w_{ij} = \text{softmax}_j(Q^l h_i^l \cdot K^l h_j^l)$$

then (17) could be written as

$$h_i^{l+1} = \sum_{j \in N(i)} w_{ij} (V^l h_j^l) = W(h_i^l)^T \sum_{j \in N(i)} (V^l h_j^l)$$

which is a special case for (15)

(e) This is because for fully-connected graphs, number of edges in the graph scales quadratically with the number of nodes. So for sentence with n words, the Transformer / GNN will compute n^2 pairs of words. Thus for a very-long-term n , the computation will be problematic.

Problem 6

Briefly summarize the key contributions, strengths and weaknesses of this paper.

Main Contribution: In this paper, Park et. al propose a multimodal explanation model that provides natural language justification that explains the content in an image. They also built a new dataset that defines visual and textual justifications of a classification decision for activity recognition tasks (ACT-X) and for visual question answering tasks (VQA-X). They suggest that training with the textual explanations yields better textual justification models and better localizes the evidence that supports the decision.

Strength: They proposed an creative idea from traditional recognition scheme by "reasoning". Their textual justifications could generate both generic and non-generic details of an image, meaning that the model has "understandings" of image contents.

Weakness: I think from the Textual Clarification evaluation table, though the scores showing improvement when there are no description or attention, the scores are mostly below 50%. Though the final Multimodal Explanations score is pretty good, I'm wondering if improving the Textual Clarification score will affect the Multimodal Explanations score.

Problem 7

What is your personal takeaway from this paper? This could be expressed either in terms of relating the approaches adopted in this paper to your traditional understanding of learning parameterized models, or potential future directions of research in the area which the authors haven't addressed, or anything else that struck you as being noteworthy.

Instead of just feeding image to a network, this paper proposed a multimodal model that allows a third-person, post-hoc rationalization type of explanation of images, akin to what one human produces when asked to explain to another human. This would also be helpful for image recognition researchers to know what's going on under the hood rather than just feed image into the networks. Future direction might be improving the Textual Clarification scores and improve the model in general.

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.

```
In [6]: # As usual, a bit of setup
from __future__ import print_function
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_g
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
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-i
%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))))
```

Install h5py

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

```
pip install h5py
```

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

```
sudo pip install h5py
```

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

In [3]: !pip install h5py

```
Collecting h5py
  Downloading h5py-2.10.0-cp37-cp37m-macosx_10_6_intel.whl (3.0 MB)
    |██████████| 3.0 MB 2.4 MB/s eta 0:00:01
Requirement already satisfied: numpy>=1.7 in /opt/anaconda3/envs/cs7643/lib/python3.7/site-packages (from h5py) (1.19.1)
Requirement already satisfied: six in /opt/anaconda3/envs/cs7643/lib/python3.7/site-packages (from h5py) (1.15.0)
Installing collected packages: h5py
Successfully installed h5py-2.10.0
```

Microsoft COCO

For this exercise we will use the 2014 release of the [Microsoft COCO dataset \(`http://mscoco.org/`\)](http://mscoco.org/), 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:

```
In [7]: # 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))

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
idxToWord <class 'list'> 1004
wordToIdx <class 'dict'> 1004
trainUrls <class 'numpy.ndarray'> (82783,) <U63
valUrls <class 'numpy.ndarray'> (40504,) <U63
```

Look at the data

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

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

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

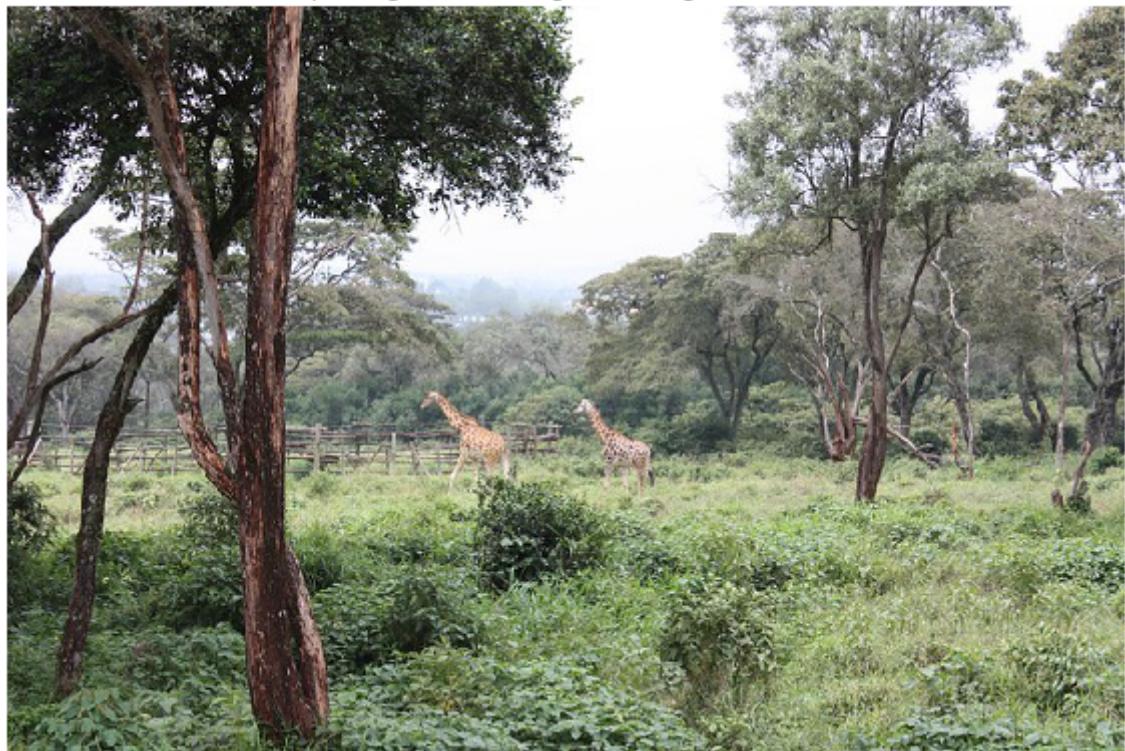
```
In [8]: # Sample a minibatch and show the images and captions
batch_size = 3

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_caption(caption, data['idx_to_word'])
    plt.title(caption_str)
    plt.show()
```

<START> the man <UNK> <UNK> the tennis ball with his racket <END>



<START> a couple of giraffe walking across a grass covered field <END>



<START> a herd of <UNK> grazing grassy rocky land <END>



Recurrent Neural Networks

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

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

Vanilla RNN: step forward

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

First implement the function `rnn_step_forward` which implements the forward pass for a single timestep of a vanilla recurrent neural network. After doing so run the following to check your implementation. You should see errors less than 1e-8.

```
In [11]: N, D, H = 3, 10, 4
```

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

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

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

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

Vanilla RNN: step backward

In the file `cs231n/rnn_layers.py` implement the `rnn_step_backward` function. After doing so run the following to numerically gradient check your implementation. You should see errors less than `1e-8`.

```
In [13]: 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
```

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 process 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 less than `1e-7`.

In [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

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, calling into the `rnn_step_backward` function that you defined above. You should see errors less than 5e-7.

```
In [20]: np.random.seed(231)
```

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

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

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

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

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

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

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

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

```
dx error:  1.5402322184213243e-09
dh0 error:  3.3824326261334578e-09
dWx error:  7.238350796069372e-09
dWh error:  1.3157659173166636e-07
db error:  1.5353591509855146e-10
```

Word embedding: forward

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

In the file `cs231n/rnn_layers.py`, implement the function `word_embedding_forward` to convert words (represented by integers) into vectors. Run the following to check your implementation. You should see error around `1e-8`.

```
In [21]: N, T, V, D = 2, 4, 5, 3

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

out, _ = word_embedding_forward(x, W)
expected_out = np.asarray([
    [[0., 0.07142857, 0.14285714],
     [0.64285714, 0.71428571, 0.78571429],
     [0.21428571, 0.28571429, 0.35714286],
     [0.42857143, 0.5, 0.57142857]],
    [[0.42857143, 0.5, 0.57142857],
     [0.21428571, 0.28571429, 0.35714286],
     [0., 0.07142857, 0.14285714],
     [0.64285714, 0.71428571, 0.78571429]]])

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

out error: 1.000000094736443e-08

Word embedding: backward

Implement the backward pass for the word embedding function in the function `word_embedding_backward`. After doing so run the following to numerically gradient check your implementation. You should see errors less than `1e-11`.

```
In [22]: 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

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 less than `1e-9`.

In [23]: `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
```

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 less than 1e-7.

```
In [24]: # 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, p):
    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 about 2.3

# 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))

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

2.3027781774290146
23.025985953127226
2.2643611790293394
dx error: 2.583585303524283e-08

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 vanilla 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 less than `1e-10`.

```
In [41]: 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 errors around $5e-6$ or less.

```
In [44]: 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], v
    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
```

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 overfit a small sample of 100 training examples. You should see losses of less than 0.1.

```
In [45]: np.random.seed(231)
```

```
small_data = load_coco_data(max_train=100)

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,
)

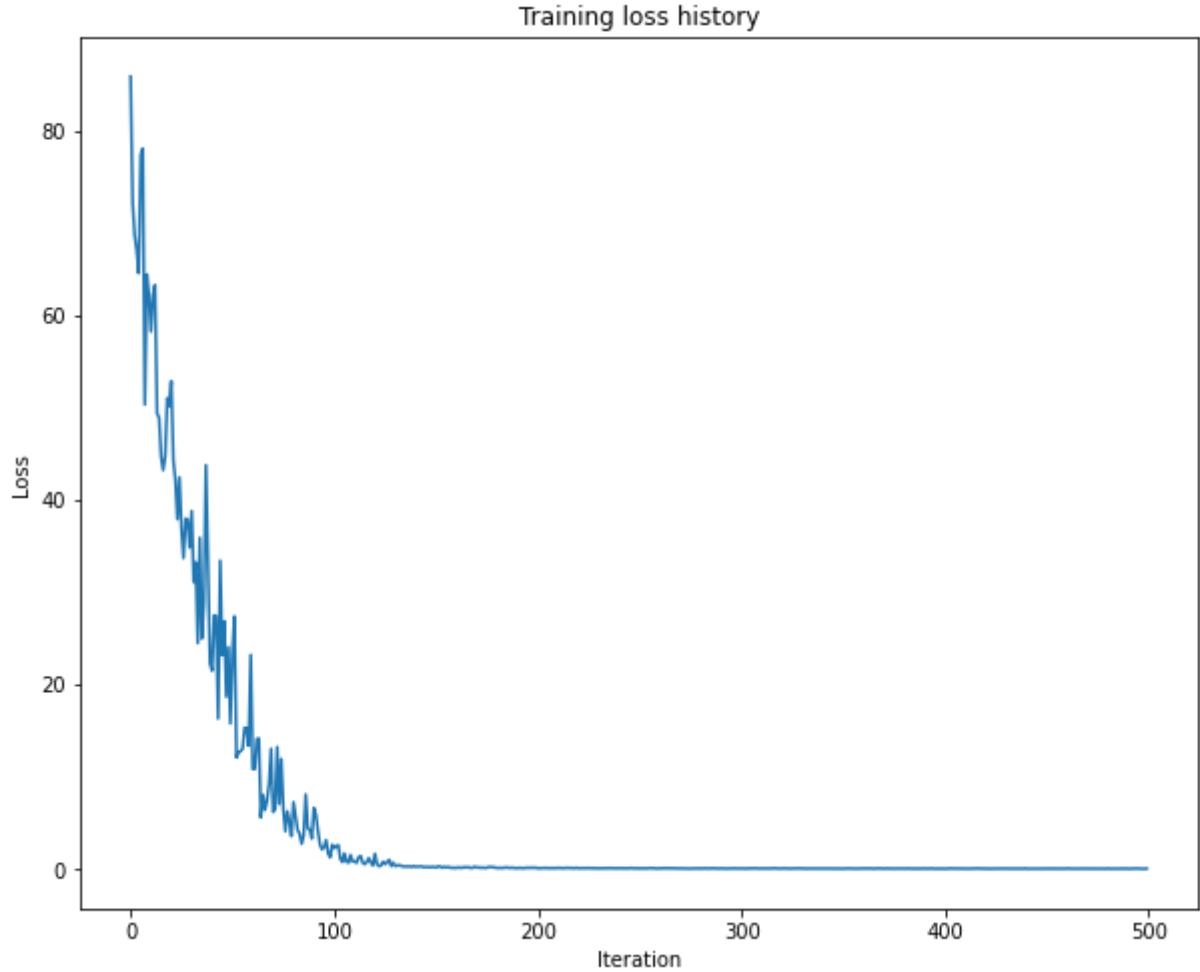
small_rnn_solver = CaptioningSolver(small_rnn_model, small_data,
    update_rule='adam',
    num_epochs=50,
    batch_size=10,
    optim_config={
        'learning_rate': 5e-3,
    },
    lr_decay=0.95,
    verbose=True, print_every=10,
)

small_rnn_solver.train()

# Plot the training losses
plt.plot(small_rnn_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```

```
(Iteration 1 / 500) loss: 85.893020
(Iteration 11 / 500) loss: 58.256967
(Iteration 21 / 500) loss: 52.886655
(Iteration 31 / 500) loss: 38.818996
(Iteration 41 / 500) loss: 21.517494
(Iteration 51 / 500) loss: 23.474711
(Iteration 61 / 500) loss: 10.860008
(Iteration 71 / 500) loss: 6.211922
(Iteration 81 / 500) loss: 7.267651
(Iteration 91 / 500) loss: 6.646090
(Iteration 101 / 500) loss: 2.280996
(Iteration 111 / 500) loss: 0.876581
(Iteration 121 / 500) loss: 1.683378
(Iteration 131 / 500) loss: 0.358072
(Iteration 141 / 500) loss: 0.287243
(Iteration 151 / 500) loss: 0.193348
(Iteration 161 / 500) loss: 0.153990
(Iteration 171 / 500) loss: 0.207464
(Iteration 181 / 500) loss: 0.137890
(Iteration 191 / 500) loss: 0.117546
(Iteration 201 / 500) loss: 0.110117
(Iteration 211 / 500) loss: 0.151787
(Iteration 221 / 500) loss: 0.110791
(Iteration 231 / 500) loss: 0.117869
```

```
(Iteration 241 / 500) loss: 0.119163
(Iteration 251 / 500) loss: 0.088646
(Iteration 261 / 500) loss: 0.110261
(Iteration 271 / 500) loss: 0.097387
(Iteration 281 / 500) loss: 0.088180
(Iteration 291 / 500) loss: 0.104136
(Iteration 301 / 500) loss: 0.094128
(Iteration 311 / 500) loss: 0.090512
(Iteration 321 / 500) loss: 0.089518
(Iteration 331 / 500) loss: 0.085879
(Iteration 341 / 500) loss: 0.092307
(Iteration 351 / 500) loss: 0.078187
(Iteration 361 / 500) loss: 0.093489
(Iteration 371 / 500) loss: 0.070030
(Iteration 381 / 500) loss: 0.077482
(Iteration 391 / 500) loss: 0.076031
(Iteration 401 / 500) loss: 0.055339
(Iteration 411 / 500) loss: 0.082718
(Iteration 421 / 500) loss: 0.081242
(Iteration 431 / 500) loss: 0.065360
(Iteration 441 / 500) loss: 0.068113
(Iteration 451 / 500) loss: 0.075713
(Iteration 461 / 500) loss: 0.091829
(Iteration 471 / 500) loss: 0.072376
(Iteration 481 / 500) loss: 0.087877
(Iteration 491 / 500) loss: 0.060424
```



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.

Note: Some of the URLs are missing and will throw an error; re-run this cell until the output is at least 2 good caption samples.

```
In [47]: 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  
        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 woman is kneeling near some large <UNK> of food <END>
GT:<START> a woman is kneeling near some large <UNK> of food <END>



train
a group of men riding in a boat across a lake <END>
GT:<START> a group of men riding in a boat across a lake <END>



val
half a <UNK> to the camera and of <UNK> <END>
GT:<START> the man in the helmet is jumping while wearing <UNK> <UNK> <UNK> <END>



val

two with the eating <UNK> <UNK> <END>

GT:<START> a little boy sitting on the stairs with a racquet <END>



In []:

Image Captioning with LSTMs

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

```
In [4]: # As usual, a bit of setup
from __future__ import print_function
import time, os, json
import numpy as np
import matplotlib.pyplot as plt
import nltk

from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_g
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
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-i
%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))))
```

The autoreload extension is already loaded. To reload it, use:
`%reload_ext autoreload`

Load MS-COCO data

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

```
In [5]: # 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))
```

```
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
```

LSTM

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

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

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

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

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

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

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

where \odot is the elementwise product of vectors.

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

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

LSTM: step forward

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

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

```
In [7]: 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
```

LSTM: step backward

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

In [9]: `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: 1.6933643922734908e-09
dWh error: 4.806248540056623e-08
db error: 1.734924139321044e-10

```

LSTM: forward

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

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

```
In [10]: N, D, H, T = 2, 5, 4, 3
x = np.linspace(-0.4, 0.6, num=N*T*D).reshape(N, T, D)
h0 = np.linspace(-0.4, 0.8, num=N*H).reshape(N, H)
Wx = np.linspace(-0.2, 0.9, num=4*D*H).reshape(D, 4 * H)
Wh = np.linspace(-0.3, 0.6, num=4*H*H).reshape(H, 4 * H)
b = np.linspace(0.2, 0.7, num=4*H)

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

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

print('h error: ', rel_error(expected_h, h))

h error:  8.610537452106624e-08
```

LSTM: backward

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

```
In [14]: 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:  6.186239731971807e-09
dh0 error:  6.791931828439997e-09
dWx error:  3.301449268876082e-09
dWh error:  1.5077886220706733e-06
db error:  7.214755671033742e-10
```

LSTM captioning model

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

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

```
In [15]: 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
```

Overfit LSTM captioning model

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

```
In [16]: np.random.seed(231)
```

```
small_data = load_coco_data(max_train=50)

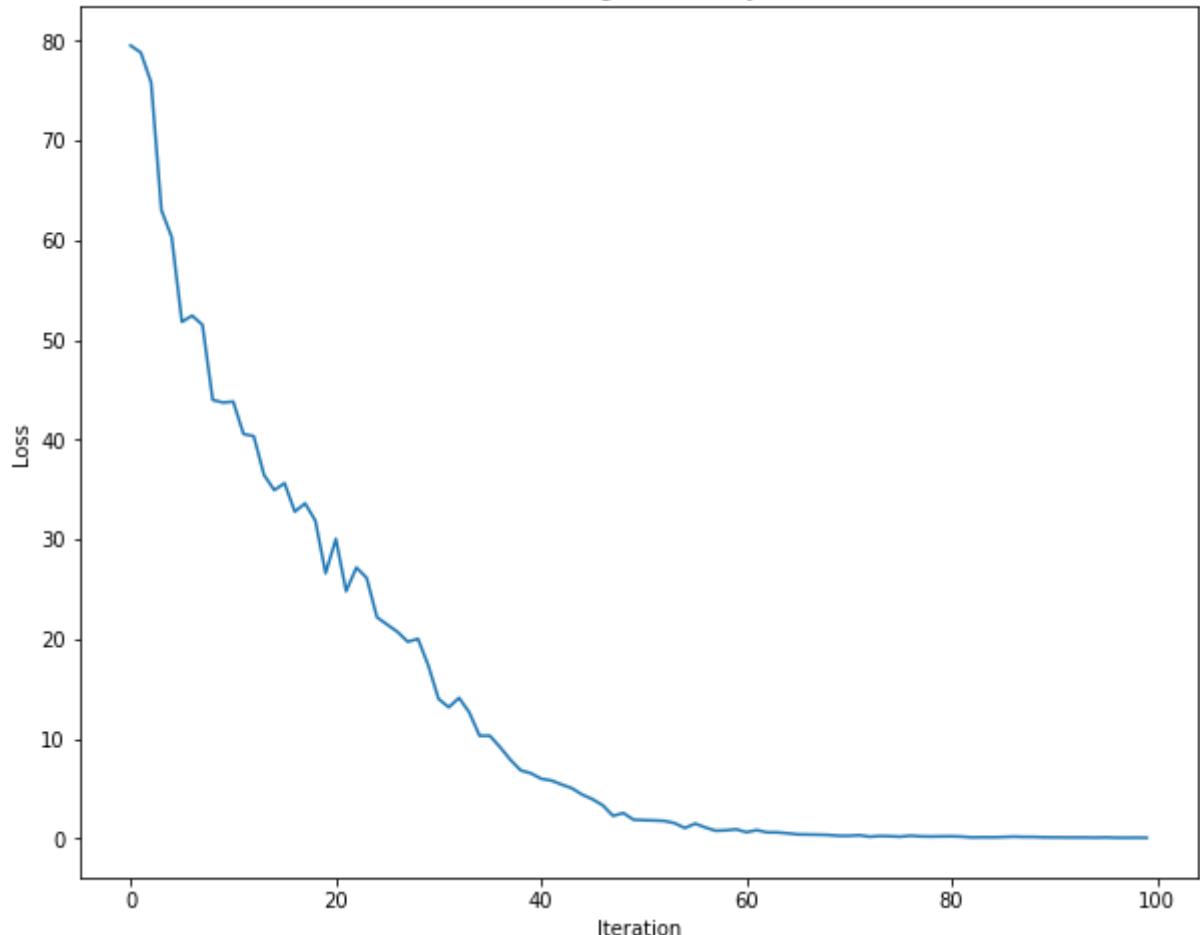
small_lstm_model = CaptioningRNN(
    cell_type='lstm',
    word_to_idx=data['word_to_idx'],
    input_dim=data['train_features'].shape[1],
    hidden_dim=512,
    wordvec_dim=256,
    dtype=np.float32,
)

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

# Plot the training losses
plt.plot(small_lstm_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```

```
(Iteration 1 / 100) loss: 79.551150
(Iteration 11 / 100) loss: 43.829099
(Iteration 21 / 100) loss: 30.062605
(Iteration 31 / 100) loss: 14.020170
(Iteration 41 / 100) loss: 6.005865
(Iteration 51 / 100) loss: 1.851858
(Iteration 61 / 100) loss: 0.643203
(Iteration 71 / 100) loss: 0.283403
(Iteration 81 / 100) loss: 0.232544
(Iteration 91 / 100) loss: 0.124545
```

Training loss history



LSTM test-time sampling

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

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

```
In [22]: for split in ['train', 'val']:  
    minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2  
    gt_captions, features, urls = minibatch  
    gt_captions = decode_captions(gt_captions, data['idx_to_word'])  
  
    sample_captions = small_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:  
        plt.imshow(image_from_url(url))  
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))  
        plt.axis('off')  
        plt.show()
```

train
many people standing near boxes of many apples <END>
GT:<START> many people standing near boxes of many apples <END>



train

a surfer rides a large wave while the sun <UNK> the <UNK> <END>
GT:<START> a surfer rides a large wave while the sun <UNK> the <UNK> <END>



val

filled five grazing grazing grazing sleeping cute cute dog standing on a the side of a <END>
GT:<START> a bowl of chicken and vegetables is shown <END>



val

an open refrigerator standing with a man on a box of suitcases <END>
GT:<START> a salad and a sandwich <UNK> to be eaten at a restaurant <END>



Train a good captioning model (extra credit for both 4803 and 7643)

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

In addition to qualitatively evaluating your model by inspecting its results, you can also quantitatively evaluate your model using the BLEU unigram precision metric. In order to achieve full credit you should train a model that achieves a BLEU unigram score of >0.25. BLEU scores range from 0 to 1; the closer to 1, the better. Here's a reference to the [paper](http://www.aclweb.org/anthology/P02-1040.pdf) (<http://www.aclweb.org/anthology/P02-1040.pdf>) that introduces BLEU if you're interested in learning more about how it works.

Feel free to use PyTorch for this section if you'd like to train faster on a GPU... though you can definitely get above 0.25 using your Numpy code. We're providing you the evaluation code that is compatible with the Numpy model as defined above... you should be able to adapt it for PyTorch if you go that route.

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

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

```
In [23]: def BLEU_score(gt_caption, sample_caption):
    """
        gt_caption: string, ground-truth caption
        sample_caption: string, your model's predicted caption
        Returns unigram BLEU score.
    """
    reference = [x for x in gt_caption.split(' ')]
        if ('<END>' not in x and '<START>' not in x and '<UNK>' no
hypothesis = [x for x in sample_caption.split(' ')]
        if ('<END>' not in x and '<START>' not in x and '<UNK>' n
BLEUscore = nltk.translate.bleu_score.sentence_bleu([reference], hypothesis)
    return BLEUscore

def evaluate_model(model):
    """
        model: CaptioningRNN model
        Prints unigram BLEU score averaged over 1000 training and val examples.
    """
    BLEUscores = {}
    for split in ['train', 'val']:
        minibatch = sample_coco_minibatch(data, split=split, batch_size=100)
        gt_captions, features, urls = minibatch
        gt_captions = decode_captions(gt_captions, data['idx_to_word'])

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

        total_score = 0.0
        for gt_caption, sample_caption, url in zip(gt_captions, sample_captions, urls):
            total_score += BLEU_score(gt_caption, sample_caption)

        BLEUscores[split] = total_score / len(sample_captions)

    for split in BLEUscores:
        print('Average BLEU score for %s: %f' % (split, BLEUscores[split]))
```

```
In [ ]: # write a description of your model here:
```

```
In [ ]: # write your code to train your model here.
# make sure to include the call to evaluate_model which prints out your hig
```

Sentence Classification with Transformers

In this exercise you will implement a [Transformer](https://arxiv.org/pdf/1706.03762.pdf) (<https://arxiv.org/pdf/1706.03762.pdf>) and use it to judge the grammaticality of English sentences.

A quick note: if you receive the following TypeError "super(type, obj): obj must be an instance or subtype of type", try restarting your kernel and re-running all cells. Once you have finished making changes to the model constructor, you can avoid this issue by commenting out all of the model instantiations after the first (e.g. lines starting with "model = ClassificationTransformer").

```
In [1]: import numpy as np
import csv
import torch

from gt_7643.transformer import ClassificationTransformer

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

The Corpus of Linguistic Acceptability (CoLA)

The Corpus of Linguistic Acceptability ([CoLA](https://nyu-mll.github.io/CoLA/) (<https://nyu-mll.github.io/CoLA/>)) in its full form consists of 10657 sentences from 23 linguistics publications, expertly annotated for acceptability (grammaticality) by their original authors. Native English speakers consistently report a sharp contrast in acceptability between pairs of sentences. Some examples include:

What did Betsy paint a picture of? (Correct)

What was a picture of painted by Betsy? (Incorrect)

You can read more info about the dataset [here](https://arxiv.org/pdf/1805.12471.pdf) (<https://arxiv.org/pdf/1805.12471.pdf>). This is a binary classification task (predict 1 for correct grammar and 0 otherwise).

Can we train a neural network to accurately predict these human acceptability judgements? In this assignment, we will implement the forward pass of the Transformer architecture discussed in class. The general intuitive notion is that we will *encode* the sequence of tokens in the sentence, and then predict a binary output based on the final state that is the output of the model.

Load the preprocessed data

We've appended a "CLS" token to the beginning of each sequence, which can be used to make predictions. The benefit of appending this token to the beginning of the sequence (rather than the end) is that we can extract it quite easily (we don't need to remove paddings and figure out the length of each individual sequence in the batch). We'll come back to this.

We've additionally already constructed a vocabulary and converted all of the strings of tokens into integers which can be used for vocabulary lookup for you. Feel free to explore the data here.

```
In [2]: train_inxs = np.load('./gt_7643/datasets/train_inxs.npy')
val_inxs = np.load('./gt_7643/datasets/val_inxs.npy')
train_labels = np.load('./gt_7643/datasets/train_labels.npy')
val_labels = np.load('./gt_7643/datasets/val_labels.npy')

# load dictionary
word_to_ix = {}
with open("./gt_7643/datasets/word_to_ix.csv", "r") as f:
    reader = csv.reader(f)
    for line in reader:
        word_to_ix[line[0]] = line[1]
print("Vocabulary Size:", len(word_to_ix))

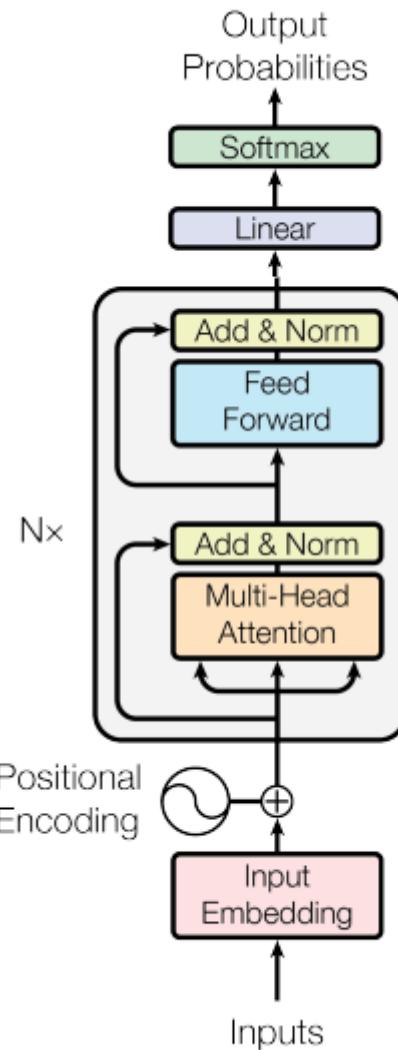
print(train_inxs.shape) # 7000 training instances, of (maximum/padded) leng
print(val_inxs.shape) # 1551 validation instances, of (maximum/padded) leng
print(train_labels.shape)
print(val_labels.shape)

# load checkers
d1 = torch.load('./gt_7643/datasets/d1.pt')
d2 = torch.load('./gt_7643/datasets/d2.pt')
d3 = torch.load('./gt_7643/datasets/d3.pt')
d4 = torch.load('./gt_7643/datasets/d4.pt')

Vocabulary Size: 1542
(7000, 43)
(1551, 43)
(7000,)
(1551,)
```

Transformers

We will be implementing a one-layer Transformer **encoder** which, similar to an RNN, can encode a sequence of inputs and produce a final output state for classification. This is the architecture:



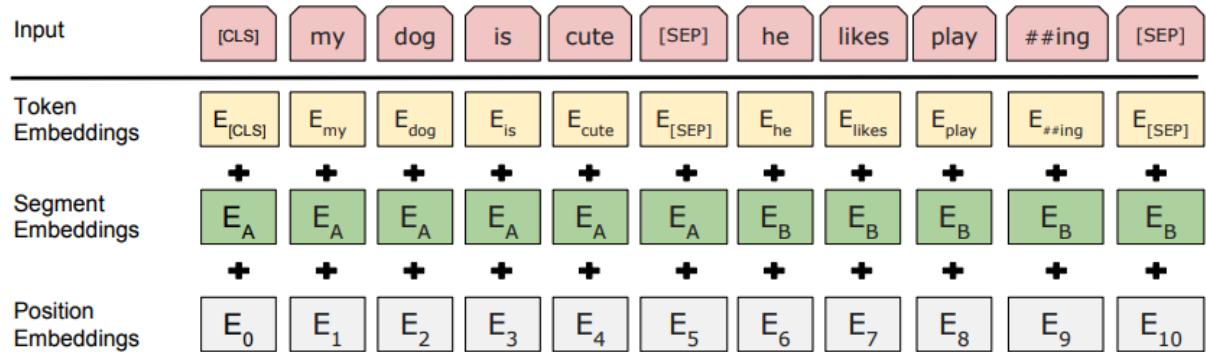
You can refer to the [original paper \(<https://arxiv.org/pdf/1706.03762.pdf>\)](https://arxiv.org/pdf/1706.03762.pdf) for more details.

Instead of using numpy for this model, we will be using Pytorch to implement the forward pass. You will not need to implement the backward pass for the various layers in this assignment.

The file `gt_7643/transformer.py` contains the model class and methods for each layer. This is where you will write your implementations.

Deliverable 1: Embeddings

We will format our input embeddings similarly to how they are constructed in [BERT \(source of figure\) \(<https://arxiv.org/pdf/1810.04805.pdf>\)](https://arxiv.org/pdf/1810.04805.pdf). Recall from lecture that unlike a RNN, a Transformer does not include any positional information about the order in which the words in the sentence occur. Because of this, we need to append a positional encoding token at each position. (We will ignore the segment embeddings and [SEP] token here, since we are only encoding one sentence at a time). We have already appended the [CLS] token for you in the previous step.



Your first task is to implement the embedding lookup, including the addition of positional encodings. Open the file `gt_7643/transformer.py` and complete all code parts for Deliverable 1.

```
In [3]: inputs = train_inxs[0:2]
inputs = torch.LongTensor(inputs)

model = ClassificationTransformer(word_to_ix, hidden_dim=128, num_heads=2,
                                  dim_v=96, dim_q=96, max_length=train_inxs

embeds = model.embed(inputs)

try:
    print("Difference:", torch.sum(torch.pairwise_distance(embeds, d1)).item)
except:
    print("NOT IMPLEMENTED")
```

Difference: 0.0017988268518820405

Deliverable 2: Multi-head Self-Attention

Attention can be computed in matrix-form using the following formula:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

We want to have multiple self-attention operations, computed in parallel. Each of these is called a *head*. We concatenate the heads and multiply them with the matrix `attention_head_projection` to produce the output of this layer.

After every multi-head self-attention and feedforward layer, there is a residual connection + layer normalization. Make sure to implement this, using the following formula:

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

Open the file `gt_7643/transformer.py` and implement the `multihead_attention` function. We have already initialized all of the layers you will need in the constructor.

```
In [4]: hidden_states = model.multi_head_attention(embeds)

try:
    print("Difference:", torch.sum(torch.pairwise_distance(hidden_states, d
except:
    print("NOT IMPLEMENTED")
```

Difference: 0.0017100314144045115

Deliverable 3: Element-Wise Feed-forward Layer

Open the file `gt_7643/transformer.py` and complete code for Deliverable 3 : the element-wise feed-forward layer consisting of two linear transformers with a ReLU layer in between.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

```
In [5]: model = ClassificationTransformer(word_to_ix, hidden_dim=128, num_heads=2,
                                         dim_v=96, dim_q=96, max_length=train_inxs

outputs = model.feedforward_layer(hidden_states)

try:
    print("Difference:", torch.sum(torch.pairwise_distance(outputs, d3)).item()
except:
    print("NOT IMPLEMENTED")
```

Difference: 0.0017168164486065507

Deliverable 4: Final Layer

Open the file `gt_7643/transformer.py` and complete code for Deliverable 4 , to produce binary classification scores for the inputs based on the output of the Transformer.

```
In [6]: model = ClassificationTransformer(word_to_ix, hidden_dim=128, num_heads=2,
                                         dim_v=96, dim_q=96, max_length=train_inxs

scores = model.final_layer(outputs)

try:
    print("Difference:", torch.sum(torch.pairwise_distance(scores, d4)).item()
except:
    print("NOT IMPLEMENTED")
```

Difference: 1.8956918665935518e-06

Deliverable 5: Putting it all together

Open the file `gt_7643/transformer.py` and complete the method `forward`, by putting together all of the methods you have developed in the right order to perform a full forward pass.

```
In [7]: inputs = train_inxs[0:2]
inputs = torch.LongTensor(inputs)

outputs = model.forward(inputs)

try:
    print("Difference:", torch.sum(torch.pairwise_distance(outputs, scores))
except:
    print("NOT IMPLEMENTED")
```

Difference: 1.999999949504854e-06

Great! We've just implemented a Transformer forward pass for text classification. One of the big perks of using PyTorch is that with a simple training loop, we can rely on automatic differentiation ([autograd \(\[https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html\]\(https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html\)\)](https://pytorch.org/tutorials/beginner/blitz/autograd_tutorial.html)) to do the work of the backward pass for us. This is not required for this assignment, but you can explore this on your own.

Make sure when you submit your PDF for this assignment to also include a copy of `transformer.py` converted to PDF as well.

```
In [ ]:
```

```
import numpy as np
```

```
import torch
```

```
from torch import nn
```

```
import random
```

Do not modify these imports.

```
class ClassificationTransformer(nn.Module):
```

```
    """
```

A single-layer Transformer which encodes a sequence of text and performs binary classification.

The model has a vocab size of V , works on sequences of length T , has an hidden dimension of H , uses word vectors also of dimension H , and operates on minibatches of size N .

```
    """
```

```
def __init__(self, word_to_ix, hidden_dim=128, num_heads=2, dim_feedforward=2048, dim_k=96, dim_v=96, dim_q=96, max_length=43):
```

```
    """
```

:param word_to_ix: dictionary mapping words to unique indices

:param hidden_dim: the dimensionality of the output embeddings that go into the final layer

:param num_heads: the number of Transformer heads to use

:param dim_feedforward: the dimension of the feedforward network model

:param dim_k: the dimensionality of the key vectors

:param dim_q: the dimensionality of the query vectors

:param dim_v: the dimensionality of the value vectors

```
    """
```

```
super(ClassificationTransformer, self).__init__()
```

```
assert hidden_dim % num_heads == 0
```

```
self.num_heads = num_heads
```

```
self.word_embedding_dim = hidden_dim
```

```
self.hidden_dim = hidden_dim
```

```
self.dim_feedforward = dim_feedforward
```

```
self.max_length = max_length
```

```
self.vocab_size = len(word_to_ix)
```

```
self.dim_k = dim_k
```

```
self.dim_v = dim_v
```

```
self.dim_q = dim_q
```

```
seed_torch(0)
```

```
#####
```

```
# Deliverable 1: Initialize what you need for the embedding lookup.      #
```

```
# You will need to use the max_length parameter above.                  #
```

```
# This should take 1-2 lines.                                         #
```

```
# Initialize the word embeddings before the positional encodings.      #
```

```
# Don't worry about sine/cosine encodings- use positional encodings. #
```

```
#####
```

```
self.token_embed = nn.Embedding(num_embeddings = self.vocab_size, embedding_dim = self.hidden_dim)
```

```
self.posit_embed = nn.Embedding(num_embeddings = self.max_length, embedding_dim = self.hidden_dim)
```

```
#####
```

```
# END OF YOUR CODE          #
```

```
#####
```

```
#####
```

```
# Deliverable 2: Initializations for multi-head self-attention.      #
```

```
# You don't need to do anything here. Do not modify this code.       #
```

```
#####
```

```

# Head #1
self.k1 = nn.Linear(self.hidden_dim, self.dim_k)
self.v1 = nn.Linear(self.hidden_dim, self.dim_v)
self.q1 = nn.Linear(self.hidden_dim, self.dim_q)

# Head #2
self.k2 = nn.Linear(self.hidden_dim, self.dim_k)
self.v2 = nn.Linear(self.hidden_dim, self.dim_v)
self.q2 = nn.Linear(self.hidden_dim, self.dim_q)

self.softmax = nn.Softmax(dim=2)
self.attention_head_projection = nn.Linear(self.dim_v * self.num_heads, self.hidden_dim)
self.norm_mh = nn.LayerNorm(self.hidden_dim)

#####
# Deliverable 3: Initialize what you need for the feed-forward layer.      #
# Don't forget the layer normalization.                                     #
#####
self.ffLayer1 = nn.Linear(self.hidden_dim, self.dim_feedforward)
self.reluLayer = nn.ReLU()
self.ffLayer2 = nn.Linear(self.dim_feedforward, self.hidden_dim)
self.ff_norm = nn.LayerNorm(self.hidden_dim)

#####
#           END OF YOUR CODE          #
#####

#####

# Deliverable 4: Initialize what you need for the final layer (1-2 lines). #
#####
self.finalLayer = nn.Linear(self.hidden_dim, 1)
self.sigmoid = nn.Sigmoid()

#####
#           END OF YOUR CODE          #
#####

def forward(self, inputs):
    """
    This function computes the full Transformer forward pass.
    Put together all of the layers you've developed in the correct order.

    :param inputs: a PyTorch tensor of shape (N,T). These are integer lookups.

    :returns: the model outputs. Should be normalized scores of shape (N,1).
    """

    outputs = None
    #####
    # Deliverable 5: Implement the full Transformer stack for the forward pass. #
    # You will need to use all of the methods you have previously defined above.#
    # You should only be calling ClassificationTransformer class methods here. #
    #####
    embeds = self.embed(inputs)
    hidden_states = self.multi_head_attention(embeds)
    FFoutputs = self.feedforward_layer(hidden_states)
    outputs = self.final_layer(FFoutputs)

    #####
    #           END OF YOUR CODE          #
    #####
    return outputs

```

```
def embed(self, inputs):
```

```

"""
:param inputs: intTensor of shape (N, T)
:returns embeddings: floatTensor of shape (N, T, H)
"""

embeddings = None
#####
# Deliverable 1: Implement the embedding lookup. #
# Note: word_to_ix has keys from 0 to self.vocab_size - 1 #
# This will take a few lines. #
#####
#N, T = inputs.shape
embeddings = torch.add(self.token_embed(inputs), self.posit_embed(torch.arange(inputs.shape[1])))

#####
# END OF YOUR CODE #
#####

return embeddings

def multi_head_attention(self, inputs):
"""
:param inputs: float32 Tensor of shape (N, T, H)
:returns outputs: float32 Tensor of shape (N, T, H)

Traditionally we'd include a padding mask here, so that pads are ignored.
This is a simplified implementation.
"""

outputs = None
#####
# Deliverable 2: Implement multi-head self-attention followed by add + norm.#
# Use the provided 'Deliverable 2' layers initialized in the constructor. #
#####

Q1 = self.q1(inputs)
K1 = self.k1(inputs)
V1 = self.v1(inputs)
#print(Q1.size())
#print(K1.size())
#print(V1.size())
#print(K1.transpose(-2, -1).shape)
att1 = torch.matmul(self.softmax(torch.matmul(Q1, K1.transpose(-2, -1)) / np.sqrt(self.dim_k)), V1)

Q2 = self.q2(inputs)
K2 = self.k2(inputs)
V2 = self.v2(inputs)
att2 = torch.matmul(self.softmax(torch.matmul(Q2, K2.transpose(-2, -1)) / np.sqrt(self.dim_k)), V2)

head = torch.cat((att1, att2), dim=2)
outputs = self.attention_head_projection(head)
outputs = self.norm_mh(torch.add(outputs, inputs))

#####
# END OF YOUR CODE #
#####

return outputs

def feedforward_layer(self, inputs):
"""
:param inputs: float32 Tensor of shape (N, T, H)
:returns outputs: float32 Tensor of shape (N, T, H)
"""

outputs = None
#####
# Deliverable 3: Implement the feedforward layer followed by add + norm. #
# Use a ReLU activation and apply the linear layers in the order you #
# initialized them. #
# This should not take more than 3-5 lines of code. #
#####

```

```
ff1_out = self.ffLayer1(inputs)
relu_out = self.reluLayer(ff1_out)
ff2_out = self.ffLayer2(relu_out)
outputs = self.ff_norm(ff2_out + inputs)

#####
#           END OF YOUR CODE           #
#####

return outputs
```

```
def final_layer(self, inputs):
    """
    :param inputs: float32 Tensor of shape (N,T,H)
    :returns outputs: float32 Tensor of shape (N,1)
    """

    outputs = None
    #####
    # Deliverable 4: Implement the final layer for the Transformer classifier. #
    # This should not take more than 2 lines of code.                         #
    #####
    final_out = self.finalLayer(inputs[:,0,:].squeeze(1))
    outputs = self.sigmoid(final_out)

    #####
    #           END OF YOUR CODE           #
    #####
    return outputs
```

```
def seed_torch(seed=0):
    random.seed(seed)
    np.random.seed(seed)
    torch.manual_seed(seed)
    torch.cuda.manual_seed(seed)
    torch.backends.cudnn.benchmark = False
    torch.backends.cudnn.deterministic = True
```