

# PyTorch

November 14, 2024

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cse493g1/assignments/assignment4/'
FOLDERNAME = 'cse493g1/assignments/assignment4'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My\ Drive/{}'.format(FOLDERNAME))

# This downloads the COCO dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cse493g1/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My\ Drive/cse493g1/assignments/assignment4/cse493g1/datasets
/content/drive/My\ Drive/cse493g1/assignments/assignment4
```

## 1 Introduction to PyTorch

You've written a lot of code in this course to provide a whole host of neural network functionality. Dropout, Batch Norm, and 2D convolutions are some of the workhorses of deep learning in computer vision. You've also worked hard to make your code efficient and vectorized.

For the last part of this assignment, though, we're going to leave behind your beautiful codebase and instead migrate to one of two popular deep learning frameworks: in this instance, PyTorch (or TensorFlow, if you choose to work with that notebook).

## 1.1 Why do we use deep learning frameworks?

- Our code will now run on GPUs! This will allow our models to train much faster. When using a framework like PyTorch or TensorFlow you can harness the power of the GPU for your own custom neural network architectures without having to write CUDA code directly (which is beyond the scope of this class).
- In this class, we want you to be ready to use one of these frameworks for your project so you can experiment more efficiently than if you were writing every feature you want to use by hand.
- We want you to stand on the shoulders of giants! TensorFlow and PyTorch are both excellent frameworks that will make your lives a lot easier, and now that you understand their guts, you are free to use them :)
- Finally, we want you to be exposed to the sort of deep learning code you might run into in academia or industry.

## 1.2 What is PyTorch?

PyTorch is a system for executing dynamic computational graphs over Tensor objects that behave similarly as numpy ndarray. It comes with a powerful automatic differentiation engine that removes the need for manual back-propagation.

## 1.3 How do I learn PyTorch?

One of the former instructors of the deep learning course at Stanford, Justin Johnson, made an excellent [tutorial](#) for PyTorch.

You can also find the detailed [API doc](#) here. If you have other questions that are not addressed by the API docs, the [PyTorch forum](#) is a much better place to ask than StackOverflow.

## 2 Table of Contents

This assignment has 5 parts. You will learn PyTorch on **three different levels of abstraction**, which will help you understand it better and prepare you for the final project.

1. Part I, Preparation: we will use CIFAR-10 dataset.
2. Part II, Barebones PyTorch: **Abstraction level 1**, we will work directly with the lowest-level PyTorch Tensors.
3. Part III, PyTorch Module API: **Abstraction level 2**, we will use `nn.Module` to define arbitrary neural network architecture.
4. Part IV, PyTorch Sequential API: **Abstraction level 3**, we will use `nn.Sequential` to define a linear feed-forward network very conveniently.
5. Part V, CIFAR-10 open-ended challenge: please implement your own network to get as high accuracy as possible on CIFAR-10. You can experiment with any layer, optimizer, hyperparameters or other advanced features.

Here is a table of comparison:

API	Flexibility	Convenience
Barebone	High	Low

API	Flexibility	Convenience
<code>nn.Module</code>	High	Medium
<code>nn.Sequential</code>	Low	High

## 3 GPU

You can manually switch to a GPU device on Colab by clicking `Runtime -> Change runtime type` and selecting GPU under `Hardware Accelerator`. You should do this before running the following cells to import packages, since the kernel gets restarted upon switching runtimes.

```
[2]: import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
from torch.utils.data import sampler

import torchvision.datasets as dset
import torchvision.transforms as T

import numpy as np

USE_GPU = True
dtype = torch.float32 # We will be using float throughout this tutorial.

if USE_GPU and torch.cuda.is_available():
    device = torch.device('cuda')
else:
    device = torch.device('cpu')

# Constant to control how frequently we print train loss.
print_every = 100
print('using device:', device)
```

using device: cuda

## 4 Part I. Preparation

Now, let's load the CIFAR-10 dataset. This might take a couple minutes the first time you do it, but the files should stay cached after that.

In previous parts of the course we had to write our own code to download the CIFAR-10 dataset, preprocess it, and iterate through it in minibatches; PyTorch provides convenient tools to automate this process for us.

```
[3]: NUM_TRAIN = 49000

# The torchvision.transforms package provides tools for preprocessing data
```

```

# and for performing data augmentation; here we set up a transform to
# preprocess the data by subtracting the mean RGB value and dividing by the
# standard deviation of each RGB value; we've hardcoded the mean and std.
transform = T.Compose([
    T.ToTensor(),
    T.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))
])

# We set up a Dataset object for each split (train / val / test); Datasets load
# training examples one at a time, so we wrap each Dataset in a DataLoader which
# iterates through the Dataset and forms minibatches. We divide the CIFAR-10
# training set into train and val sets by passing a Sampler object to the
# DataLoader telling how it should sample from the underlying Dataset.
cifar10_train = dset.CIFAR10('./cse493g1/datasets', train=True, download=True,
                             transform=transform)
loader_train = DataLoader(cifar10_train, batch_size=64,
                          sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN)))

cifar10_val = dset.CIFAR10('./cse493g1/datasets', train=True, download=True,
                           transform=transform)
loader_val = DataLoader(cifar10_val, batch_size=64,
                       sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN, ↴50000)))

cifar10_test = dset.CIFAR10('./cse493g1/datasets', train=False, download=True,
                            transform=transform)
loader_test = DataLoader(cifar10_test, batch_size=64)

```

Files already downloaded and verified

Files already downloaded and verified

Files already downloaded and verified

## 5 Part II. Barebones PyTorch

PyTorch ships with high-level APIs to help us define model architectures conveniently, which we will cover in Part II of this tutorial. In this section, we will start with the barebone PyTorch elements to understand the autograd engine better. After this exercise, you will come to appreciate the high-level model API more.

We will start with a simple fully-connected ReLU network with two hidden layers and no biases for CIFAR classification. This implementation computes the forward pass using operations on PyTorch Tensors, and uses PyTorch autograd to compute gradients. It is important that you understand every line, because you will write a harder version after the example.

When we create a PyTorch Tensor with `requires_grad=True`, then operations involving that Tensor will not just compute values; they will also build up a computational graph in the background, allowing us to easily backpropagate through the graph to compute gradients of some Tensors with respect to a downstream loss. Concretely if `x` is a Tensor with `x.requires_grad == True` then

after backpropagation `x.grad` will be another Tensor holding the gradient of `x` with respect to the scalar loss at the end.

### 5.0.1 PyTorch Tensors: Flatten Function

A PyTorch Tensor is conceptionally similar to a numpy array: it is an n-dimensional grid of numbers, and like numpy PyTorch provides many functions to efficiently operate on Tensors. As a simple example, we provide a `flatten` function below which reshapes image data for use in a fully-connected neural network.

Recall that image data is typically stored in a Tensor of shape  $N \times C \times H \times W$ , where:

- $N$  is the number of datapoints
- $C$  is the number of channels
- $H$  is the height of the intermediate feature map in pixels
- $W$  is the height of the intermediate feature map in pixels

This is the right way to represent the data when we are doing something like a 2D convolution, that needs spatial understanding of where the intermediate features are relative to each other. When we use fully connected affine layers to process the image, however, we want each datapoint to be represented by a single vector – it's no longer useful to segregate the different channels, rows, and columns of the data. So, we use a “flatten” operation to collapse the  $C \times H \times W$  values per representation into a single long vector. The flatten function below first reads in the  $N$ ,  $C$ ,  $H$ , and  $W$  values from a given batch of data, and then returns a “view” of that data. “View” is analogous to numpy’s “reshape” method: it reshapes `x`’s dimensions to be  $N \times ??$ , where  $??$  is allowed to be anything (in this case, it will be  $C \times H \times W$ , but we don’t need to specify that explicitly).

```
[4]: def flatten(x):
    N = x.shape[0] # read in N, C, H, W
    return x.view(N, -1) # "flatten" the C * H * W values into a single vector
    ↪per image

def test_flatten():
    x = torch.arange(12).view(2, 1, 3, 2)
    print('Before flattening: ', x)
    print('After flattening: ', flatten(x))

test_flatten()
```

```
Before flattening:  tensor([[[[ 0,  1],
                             [ 2,  3],
                             [ 4,  5]]],

                             [[[ 6,  7],
                               [ 8,  9],
                               [10, 11]]]])
After flattening:  tensor([[ 0,  1,  2,  3,  4,  5],
                           [ 6,  7,  8,  9, 10, 11]])
```

## 5.0.2 Barebones PyTorch: Two-Layer Network

Here we define a function `two_layer_fc` which performs the forward pass of a two-layer fully-connected ReLU network on a batch of image data. After defining the forward pass we check that it doesn't crash and that it produces outputs of the right shape by running zeros through the network.

You don't have to write any code here, but it's important that you read and understand the implementation.

```
[5]: import torch.nn.functional as F # useful stateless functions

def two_layer_fc(x, params):
    """
    A fully-connected neural networks; the architecture is:
    NN is fully connected -> ReLU -> fully connected layer.
    Note that this function only defines the forward pass;
    PyTorch will take care of the backward pass for us.

    The input to the network will be a minibatch of data, of shape
    (N, d1, ..., dM) where d1 * ... * dM = D. The hidden layer will have H_
    ↵units,
    and the output layer will produce scores for C classes.

    Inputs:
    - x: A PyTorch Tensor of shape (N, d1, ..., dM) giving a minibatch of
      input data.
    - params: A list [w1, w2] of PyTorch Tensors giving weights for the network;
      w1 has shape (D, H) and w2 has shape (H, C).

    Returns:
    - scores: A PyTorch Tensor of shape (N, C) giving classification scores for
      the input data x.
    """
    # first we flatten the image
    x = flatten(x) # shape: [batch_size, C x H x W]

    w1, w2 = params

    # Forward pass: compute predicted y using operations on Tensors. Since w1_
    ↵and
    # w2 have requires_grad=True, operations involving these Tensors will cause
    # PyTorch to build a computational graph, allowing automatic computation of
    # gradients. Since we are no longer implementing the backward pass by hand_
    ↵we
    # don't need to keep references to intermediate values.
    # you can also use `x.clamp(min=0)`_, equivalent to F.relu()
    x = F.relu(x.mm(w1))
```

```

x = x.mm(w2)
return x

def two_layer_fc_test():
    hidden_layer_size = 42
    x = torch.zeros((64, 50), dtype=dtype) # minibatch size 64, feature dimension 50
    w1 = torch.zeros((50, hidden_layer_size), dtype=dtype)
    w2 = torch.zeros((hidden_layer_size, 10), dtype=dtype)
    scores = two_layer_fc(x, [w1, w2])
    print(scores.size()) # you should see [64, 10]

two_layer_fc_test()

```

```
torch.Size([64, 10])
```

### 5.0.3 Barebones PyTorch: Three-Layer ConvNet

Here you will complete the implementation of the function `three_layer_convnet`, which will perform the forward pass of a three-layer convolutional network. Like above, we can immediately test our implementation by passing zeros through the network. The network should have the following architecture:

1. A convolutional layer (with bias) with `channel_1` filters, each with shape  $KW1 \times KH1$ , and zero-padding of two
2. ReLU nonlinearity
3. A convolutional layer (with bias) with `channel_2` filters, each with shape  $KW2 \times KH2$ , and zero-padding of one
4. ReLU nonlinearity
5. Fully-connected layer with bias, producing scores for  $C$  classes.

Note that we have **no softmax activation** here after our fully-connected layer: this is because PyTorch's cross entropy loss performs a softmax activation for you, and by bundling that step in makes computation more efficient.

**HINT:** For convolutions: <http://pytorch.org/docs/stable/nn.html#torch.nn.functional.conv2d>; pay attention to the shapes of convolutional filters!

```
[6]: def three_layer_convnet(x, params):
    """
    Performs the forward pass of a three-layer convolutional network with the
    architecture defined above.

    Inputs:
    - x: A PyTorch Tensor of shape (N, 3, H, W) giving a minibatch of images
    - params: A list of PyTorch Tensors giving the weights and biases for the
              network; should contain the following:
    - conv_w1: PyTorch Tensor of shape (channel_1, 3, KH1, KW1) giving weights
    - conv_b1: PyTorch Tensor of shape (channel_1,) giving biases for the
              first convolutional layer
    - conv_w2: PyTorch Tensor of shape (channel_2, channel_1, KH2, KW2) giving
              weights for the second convolutional layer
    - conv_b2: PyTorch Tensor of shape (channel_2,) giving biases for the
              second convolutional layer
    - fc_w: PyTorch Tensor of shape (D, C) giving weights for the fully-
           connected layer
    - fc_b: PyTorch Tensor of shape (C,) giving biases for the fully-
           connected layer
    """
    # ... forward pass code ...

```

```

    for the first convolutional layer
    - conv_b1: PyTorch Tensor of shape (channel_1,) giving biases for the
    ↵first
        convolutional layer
    - conv_w2: PyTorch Tensor of shape (channel_2, channel_1, KH2, KW2) giving
        weights for the second convolutional layer
    - conv_b2: PyTorch Tensor of shape (channel_2,) giving biases for the
    ↵second
        convolutional layer
    - fc_w: PyTorch Tensor giving weights for the fully-connected layer. Can
    ↵you
        figure out what the shape should be?
    - fc_b: PyTorch Tensor giving biases for the fully-connected layer. Can
    ↵you
        figure out what the shape should be?

>Returns:
- scores: PyTorch Tensor of shape (N, C) giving classification scores for x
"""
conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b = params
scores = None

# TODO: Implement the forward pass for the three-layer ConvNet.
#
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

out = F.conv2d(x, conv_w1, bias=conv_b1, stride=1, padding=2)
out = F.relu(out)
out = F.conv2d(out, conv_w2, bias=conv_b2, stride=1, padding=1)
out = F.relu(out)
scores = out.view(out.size(0), -1).mm(fc_w) + fc_b

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

```

After defining the forward pass of the ConvNet above, run the following cell to test your implementation.

When you run this function, scores should have shape (64, 10).

```
[7]: def three_layer_convnet_test():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image
    ↪size [3, 32, 32]

    conv_w1 = torch.zeros((6, 3, 5, 5), dtype=dtype) # [out_channel, ↪
    ↪in_channel, kernel_H, kernel_W]
    conv_b1 = torch.zeros((6,)) # out_channel
    conv_w2 = torch.zeros((9, 6, 3, 3), dtype=dtype) # [out_channel, ↪
    ↪in_channel, kernel_H, kernel_W]
    conv_b2 = torch.zeros((9,)) # out_channel

    # you must calculate the shape of the tensor after two conv layers, before
    ↪the fully-connected layer
    fc_w = torch.zeros((9 * 32 * 32, 10))
    fc_b = torch.zeros(10)

    scores = three_layer_convnet(x, [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, ↪
    ↪fc_b])
    print(scores.size()) # you should see [64, 10]
three_layer_convnet_test()
```

```
torch.Size([64, 10])
```

#### 5.0.4 Barebones PyTorch: Initialization

Let's write a couple utility methods to initialize the weight matrices for our models.

- `random_weight(shape)` initializes a weight tensor with the Kaiming normalization method.
- `zero_weight(shape)` initializes a weight tensor with all zeros. Useful for instantiating bias parameters.

The `random_weight` function uses the Kaiming normal initialization method, described in:

He et al, *Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification*, ICCV 2015, <https://arxiv.org/abs/1502.01852>

```
[8]: def random_weight(shape):
    """
    Create random Tensors for weights; setting requires_grad=True means that we
    want to compute gradients for these Tensors during the backward pass.
    We use Kaiming normalization: sqrt(2 / fan_in)
    """
    if len(shape) == 2: # FC weight
        fan_in = shape[0]
    else:
        fan_in = np.prod(shape[1:]) # conv weight [out_channel, in_channel, kH, ↪
        ↪kW]
```

```

# randn is standard normal distribution generator.
w = torch.randn(shape, device=device, dtype=dtype) * np.sqrt(2. / fan_in)
w.requires_grad = True
return w

def zero_weight(shape):
    return torch.zeros(shape, device=device, dtype=dtype, requires_grad=True)

# create a weight of shape [3 x 5]
# you should see the type `torch.cuda.FloatTensor` if you use GPU.
# Otherwise it should be `torch.FloatTensor`
random_weight((3, 5))

```

[8]:

```
tensor([[ 1.0873,  0.8092,  0.0727, -0.5770, -1.2542],
       [-0.8713,  0.3954,  0.5876,  0.1812,  0.3126],
       [-0.5077, -0.4060,  0.3706,  0.6165,  1.2231]], device='cuda:0',
      requires_grad=True)
```

## 5.0.5 Barebones PyTorch: Check Accuracy

When training the model we will use the following function to check the accuracy of our model on the training or validation sets.

When checking accuracy we don't need to compute any gradients; as a result we don't need PyTorch to build a computational graph for us when we compute scores. To prevent a graph from being built we scope our computation under a `torch.no_grad()` context manager.

[9]:

```

def check_accuracy_part2(loader, model_fn, params):
    """
    Check the accuracy of a classification model.

    Inputs:
    - loader: A DataLoader for the data split we want to check
    - model_fn: A function that performs the forward pass of the model,
        with the signature scores = model_fn(x, params)
    - params: List of PyTorch Tensors giving parameters of the model

    Returns: Nothing, but prints the accuracy of the model
    """
    split = 'val' if loader.dataset.train else 'test'
    print('Checking accuracy on the %s set' % split)
    num_correct, num_samples = 0, 0
    with torch.no_grad():
        for x, y in loader:
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.int64)
            scores = model_fn(x, params)
            _, preds = scores.max(1)

```

```

        num_correct += (preds == y).sum()
        num_samples += preds.size(0)
        acc = float(num_correct) / num_samples
        print('Got %d / %d correct (%.2f%%)' % (num_correct, num_samples, 100 * acc))
    
```

### 5.0.6 BareBones PyTorch: Training Loop

We can now set up a basic training loop to train our network. We will train the model using stochastic gradient descent without momentum. We will use `torch.functional.cross_entropy` to compute the loss; you can [read about it here](#).

The training loop takes as input the neural network function, a list of initialized parameters ([`w1`, `w2`] in our example), and learning rate.

```
[10]: def train_part2(model_fn, params, learning_rate):
    """
    Train a model on CIFAR-10.

    Inputs:
    - model_fn: A Python function that performs the forward pass of the model.
      It should have the signature scores = model_fn(x, params) where x is a
      PyTorch Tensor of image data, params is a list of PyTorch Tensors giving
      model weights, and scores is a PyTorch Tensor of shape (N, C) giving
      scores for the elements in x.
    - params: List of PyTorch Tensors giving weights for the model
    - learning_rate: Python scalar giving the learning rate to use for SGD

    Returns: Nothing
    """
    for t, (x, y) in enumerate(loader_train):
        # Move the data to the proper device (GPU or CPU)
        x = x.to(device=device, dtype=dtype)
        y = y.to(device=device, dtype=torch.long)

        # Forward pass: compute scores and loss
        scores = model_fn(x, params)
        loss = F.cross_entropy(scores, y)

        # Backward pass: PyTorch figures out which Tensors in the computational
        # graph has requires_grad=True and uses backpropagation to compute the
        # gradient of the loss with respect to these Tensors, and stores the
        # gradients in the .grad attribute of each Tensor.
        loss.backward()

        # Update parameters. We don't want to backpropagate through the
        # parameter updates, so we scope the updates under a torch.no_grad()
        # context manager to prevent a computational graph from being built.
    
```

```

    with torch.no_grad():
        for w in params:
            w -= learning_rate * w.grad

        # Manually zero the gradients after running the backward pass
        w.grad.zero_()

    if t % print_every == 0:
        print('Iteration %d, loss = %.4f' % (t, loss.item()))
        check_accuracy_part2(loader_val, model_fn, params)
        print()

```

### 5.0.7 BareBones PyTorch: Train a Two-Layer Network

Now we are ready to run the training loop. We need to explicitly allocate tensors for the fully connected weights,  $w_1$  and  $w_2$ .

Each minibatch of CIFAR has 64 examples, so the tensor shape is [64, 3, 32, 32].

After flattening,  $x$  shape should be [64, 3 \* 32 \* 32]. This will be the size of the first dimension of  $w_1$ . The second dimension of  $w_1$  is the hidden layer size, which will also be the first dimension of  $w_2$ .

Finally, the output of the network is a 10-dimensional vector that represents the probability distribution over 10 classes.

You don't need to tune any hyperparameters but you should see accuracies above 40% after training for one epoch.

```
[11]: hidden_layer_size = 4000
       learning_rate = 1e-2

       w1 = random_weight((3 * 32 * 32, hidden_layer_size))
       w2 = random_weight((hidden_layer_size, 10))

       train_part2(two_layer_fc, [w1, w2], learning_rate)
```

Iteration 0, loss = 3.6025  
 Checking accuracy on the val set  
 Got 139 / 1000 correct (13.90%)

Iteration 100, loss = 2.1441  
 Checking accuracy on the val set  
 Got 304 / 1000 correct (30.40%)

Iteration 200, loss = 2.2685  
 Checking accuracy on the val set  
 Got 365 / 1000 correct (36.50%)

Iteration 300, loss = 1.8312

```
Checking accuracy on the val set
Got 410 / 1000 correct (41.00%)
```

```
Iteration 400, loss = 1.8077
Checking accuracy on the val set
Got 439 / 1000 correct (43.90%)
```

```
Iteration 500, loss = 1.9092
Checking accuracy on the val set
Got 427 / 1000 correct (42.70%)
```

```
Iteration 600, loss = 1.5986
Checking accuracy on the val set
Got 434 / 1000 correct (43.40%)
```

```
Iteration 700, loss = 1.8117
Checking accuracy on the val set
Got 440 / 1000 correct (44.00%)
```

### 5.0.8 BareBones PyTorch: Training a ConvNet

In the below you should use the functions defined above to train a three-layer convolutional network on CIFAR. The network should have the following architecture:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You should initialize your weight matrices using the `random_weight` function defined above, and you should initialize your bias vectors using the `zero_weight` function above.

You don't need to tune any hyperparameters, but if everything works correctly you should achieve an accuracy above 42% after one epoch.

```
[12]: learning_rate = 3e-3

channel_1 = 32
channel_2 = 16

conv_w1 = None
conv_b1 = None
conv_w2 = None
conv_b2 = None
fc_w = None
fc_b = None
```

```

#####
# TODO: Initialize the parameters of a three-layer ConvNet.          #
##### *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#
def convNet(x, params) :
    conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b = params
    out = F.conv2d(x, conv_w1, bias=conv_b1, stride=1, padding=2)
    out = F.relu(out)
    out = F.conv2d(out, conv_w2, bias=conv_b2, stride=1, padding=1)
    out = F.relu(out)
    scores = out.view(out.size(0), -1).mm(fc_w) + fc_b
    return scores

conv_w1 = random_weight((channel_1, 3, 5, 5))
conv_b1 = zero_weight((channel_1,))
conv_w2 = random_weight((channel_2, channel_1, 3, 3))
conv_b2 = zero_weight((channel_2,))
fc_w = random_weight((channel_2 * 32 * 32, 10))
fc_b = zero_weight((10,))
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#           END OF YOUR CODE                                     #
#####

params = [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b]
train_part2(three_layer_convnet, params, learning_rate)

```

Iteration 0, loss = 2.7841  
 Checking accuracy on the val set  
 Got 138 / 1000 correct (13.80%)

Iteration 100, loss = 1.7803  
 Checking accuracy on the val set  
 Got 374 / 1000 correct (37.40%)

Iteration 200, loss = 1.7180  
 Checking accuracy on the val set  
 Got 397 / 1000 correct (39.70%)

Iteration 300, loss = 1.6739  
 Checking accuracy on the val set  
 Got 429 / 1000 correct (42.90%)

Iteration 400, loss = 1.6377  
 Checking accuracy on the val set  
 Got 431 / 1000 correct (43.10%)

```
Iteration 500, loss = 1.6099
Checking accuracy on the val set
Got 450 / 1000 correct (45.00%)
```

```
Iteration 600, loss = 1.5242
Checking accuracy on the val set
Got 480 / 1000 correct (48.00%)
```

```
Iteration 700, loss = 1.5630
Checking accuracy on the val set
Got 471 / 1000 correct (47.10%)
```

## 6 Part III. PyTorch Module API

Barebone PyTorch requires that we track all the parameter tensors by hand. This is fine for small networks with a few tensors, but it would be extremely inconvenient and error-prone to track tens or hundreds of tensors in larger networks.

PyTorch provides the `nn.Module` API for you to define arbitrary network architectures, while tracking every learnable parameters for you. In Part II, we implemented SGD ourselves. PyTorch also provides the `torch.optim` package that implements all the common optimizers, such as RMSProp, Adagrad, and Adam. It even supports approximate second-order methods like L-BFGS! You can refer to the [doc](#) for the exact specifications of each optimizer.

To use the Module API, follow the steps below:

1. Subclass `nn.Module`. Give your network class an intuitive name like `TwoLayerFC`.
2. In the constructor `__init__()`, define all the layers you need as class attributes. Layer objects like `nn.Linear` and `nn.Conv2d` are themselves `nn.Module` subclasses and contain learnable parameters, so that you don't have to instantiate the raw tensors yourself. `nn.Module` will track these internal parameters for you. Refer to the [doc](#) to learn more about the dozens of builtin layers. **Warning:** don't forget to call the `super().__init__()` first!
3. In the `forward()` method, define the *connectivity* of your network. You should use the attributes defined in `__init__` as function calls that take tensor as input and output the "transformed" tensor. Do *not* create any new layers with learnable parameters in `forward()`! All of them must be declared upfront in `__init__`.

After you define your Module subclass, you can instantiate it as an object and call it just like the NN forward function in part II.

### 6.0.1 Module API: Two-Layer Network

Here is a concrete example of a 2-layer fully connected network:

```
[13]: class TwoLayerFC(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super().__init__()
```

```

# assign layer objects to class attributes
self.fc1 = nn.Linear(input_size, hidden_size)
# nn.init package contains convenient initialization methods
# http://pytorch.org/docs/master/nn.html#torch-nn-init
nn.init.kaiming_normal_(self.fc1.weight)
self.fc2 = nn.Linear(hidden_size, num_classes)
nn.init.kaiming_normal_(self.fc2.weight)

def forward(self, x):
    # forward always defines connectivity
    x = flatten(x)
    scores = self.fc2(F.relu(self.fc1(x)))
    return scores

def test_TwoLayerFC():
    input_size = 50
    x = torch.zeros((64, input_size), dtype=dtype)  # minibatch size 64, ↴
    # feature dimension 50
    model = TwoLayerFC(input_size, 42, 10)
    scores = model(x)
    print(scores.size())  # you should see [64, 10]
test_TwoLayerFC()

```

`torch.Size([64, 10])`

## 6.0.2 Module API: Three-Layer ConvNet

It's your turn to implement a 3-layer ConvNet followed by a fully connected layer. The network architecture should be the same as in Part II:

1. Convolutional layer with `channel_1` 5x5 filters with zero-padding of 2
2. ReLU
3. Convolutional layer with `channel_2` 3x3 filters with zero-padding of 1
4. ReLU
5. Fully-connected layer to `num_classes` classes

You should initialize the weight matrices of the model using the Kaiming normal initialization method.

**HINT:** <http://pytorch.org/docs/stable/nn.html#conv2d>

After you implement the three-layer ConvNet, the `test_ThreeLayerConvNet` function will run your implementation; it should print `(64, 10)` for the shape of the output scores.

```
[14]: class ThreeLayerConvNet(nn.Module):
    def __init__(self, in_channel, channel_1, channel_2, num_classes):
        super().__init__()
        #####
        # TODO: Set up the layers you need for a three-layer ConvNet with the #
        # architecture defined above. #####

```

```

#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
self.convnn1 = nn.Conv2d(in_channel, channel_1, kernel_size=5, padding=2, bias=True)
nn.init.kaiming_normal_(self.convnn1.weight)
self.convnn2 = nn.Conv2d(channel_1, channel_2, kernel_size=3, padding=1, bias=True)
nn.init.kaiming_normal_(self.convnn2.weight)
self.fc = nn.Linear(channel_2 * 32 * 32, num_classes)
nn.init.kaiming_normal_(self.fc.weight)

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

def forward(self, x):
    scores = None
    #####
    # TODO: Implement the forward function for a 3-layer ConvNet. you
    # should use the layers you defined in __init__ and specify the
    # connectivity of those layers in forward()
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    scores = F.relu(self.convnn2(F.relu(self.convnn1(x))))
    scores = scores.view(scores.size(0), -1)
    scores = self.fc(scores)

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


def test_ThreeLayerConvNet():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype)  # minibatch size 64, image
    #size [3, 32, 32]
    model = ThreeLayerConvNet(in_channel=3, channel_1=12, channel_2=8, num_classes=10)
    scores = model(x)
    print(scores.size())  # you should see [64, 10]
test_ThreeLayerConvNet()

torch.Size([64, 10])

```

### 6.0.3 Module API: Check Accuracy

Given the validation or test set, we can check the classification accuracy of a neural network.

This version is slightly different from the one in part II. You don't manually pass in the parameters anymore.

```
[15]: def check_accuracy_part34(loader, model):
    if loader.dataset.train:
        print('Checking accuracy on validation set')
    else:
        print('Checking accuracy on test set')
    num_correct = 0
    num_samples = 0
    model.eval() # set model to evaluation mode
    with torch.no_grad():
        for x, y in loader:
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.long)
            scores = model(x)
            _, preds = scores.max(1)
            num_correct += (preds == y).sum()
            num_samples += preds.size(0)
        acc = float(num_correct) / num_samples
        print('Got %d / %d correct (%.2f)' % (num_correct, num_samples, 100 * acc))
```

### 6.0.4 Module API: Training Loop

We also use a slightly different training loop. Rather than updating the values of the weights ourselves, we use an Optimizer object from the `torch.optim` package, which abstract the notion of an optimization algorithm and provides implementations of most of the algorithms commonly used to optimize neural networks.

```
[16]: def train_part34(model, optimizer, epochs=1):
    """
    Train a model on CIFAR-10 using the PyTorch Module API.

    Inputs:
    - model: A PyTorch Module giving the model to train.
    - optimizer: An Optimizer object we will use to train the model
    - epochs: (Optional) A Python integer giving the number of epochs to train
    """
    for e in range(epochs):
        for t, (x, y) in enumerate(loader_train):
```

```

model.train() # put model to training mode
x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
y = y.to(device=device, dtype=torch.long)

scores = model(x)
loss = F.cross_entropy(scores, y)

# Zero out all of the gradients for the variables which the
# optimizer
# will update.
optimizer.zero_grad()

# This is the backwards pass: compute the gradient of the loss with
# respect to each parameter of the model.
loss.backward()

# Actually update the parameters of the model using the gradients
# computed by the backwards pass.
optimizer.step()

if t % print_every == 0:
    print('Iteration %d, loss = %.4f' % (t, loss.item()))
    check_accuracy_part34(loader_val, model)
    print()

```

### 6.0.5 Module API: Train a Two-Layer Network

Now we are ready to run the training loop. In contrast to part II, we don't explicitly allocate parameter tensors anymore.

Simply pass the input size, hidden layer size, and number of classes (i.e. output size) to the constructor of `TwoLayerFC`.

You also need to define an optimizer that tracks all the learnable parameters inside `TwoLayerFC`.

You don't need to tune any hyperparameters, but you should see model accuracies above 40% after training for one epoch.

```
[17]: hidden_layer_size = 4000
learning_rate = 1e-2
model = TwoLayerFC(3 * 32 * 32, hidden_layer_size, 10)
optimizer = optim.SGD(model.parameters(), lr=learning_rate)

train_part34(model, optimizer)
```

```
Iteration 0, loss = 3.2298
Checking accuracy on validation set
Got 159 / 1000 correct (15.90)
```

```
Iteration 100, loss = 2.1953
Checking accuracy on validation set
Got 360 / 1000 correct (36.00)
```

```
Iteration 200, loss = 1.8267
Checking accuracy on validation set
Got 363 / 1000 correct (36.30)
```

```
Iteration 300, loss = 2.2672
Checking accuracy on validation set
Got 370 / 1000 correct (37.00)
```

```
Iteration 400, loss = 1.8063
Checking accuracy on validation set
Got 413 / 1000 correct (41.30)
```

```
Iteration 500, loss = 2.1002
Checking accuracy on validation set
Got 386 / 1000 correct (38.60)
```

```
Iteration 600, loss = 1.7068
Checking accuracy on validation set
Got 427 / 1000 correct (42.70)
```

```
Iteration 700, loss = 1.7982
Checking accuracy on validation set
Got 428 / 1000 correct (42.80)
```

## 6.0.6 Module API: Train a Three-Layer ConvNet

You should now use the Module API to train a three-layer ConvNet on CIFAR. This should look very similar to training the two-layer network! You don't need to tune any hyperparameters, but you should achieve above 45% after training for one epoch.

You should train the model using stochastic gradient descent without momentum.

```
[18]: learning_rate = 3e-3
channel_1 = 32
channel_2 = 16

model = None
optimizer = None
#####
# TODO: Instantiate your ThreeLayerConvNet model and a corresponding optimizer #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
```

```

model = ThreeLayerConvNet(in_channel=3, channel_1=channel_1,
                           ↪channel_2=channel_2, num_classes=10)
optimizer = optim.SGD(model.parameters(), lr=learning_rate)

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

train_part34(model, optimizer)

```

Iteration 0, loss = 3.4835  
 Checking accuracy on validation set  
 Got 121 / 1000 correct (12.10)

Iteration 100, loss = 1.9267  
 Checking accuracy on validation set  
 Got 319 / 1000 correct (31.90)

Iteration 200, loss = 1.5752  
 Checking accuracy on validation set  
 Got 379 / 1000 correct (37.90)

Iteration 300, loss = 1.7198  
 Checking accuracy on validation set  
 Got 418 / 1000 correct (41.80)

Iteration 400, loss = 1.7440  
 Checking accuracy on validation set  
 Got 457 / 1000 correct (45.70)

Iteration 500, loss = 1.5609  
 Checking accuracy on validation set  
 Got 467 / 1000 correct (46.70)

Iteration 600, loss = 1.4837  
 Checking accuracy on validation set  
 Got 485 / 1000 correct (48.50)

Iteration 700, loss = 1.5202  
 Checking accuracy on validation set  
 Got 487 / 1000 correct (48.70)

## 7 Part IV. PyTorch Sequential API

Part III introduced the PyTorch Module API, which allows you to define arbitrary learnable layers and their connectivity.

For simple models like a stack of feed forward layers, you still need to go through 3 steps: subclass `nn.Module`, assign layers to class attributes in `__init__`, and call each layer one by one in `forward()`. Is there a more convenient way?

Fortunately, PyTorch provides a container Module called `nn.Sequential`, which merges the above steps into one. It is not as flexible as `nn.Module`, because you cannot specify more complex topology than a feed-forward stack, but it's good enough for many use cases.

### 7.0.1 Sequential API: Two-Layer Network

Let's see how to rewrite our two-layer fully connected network example with `nn.Sequential`, and train it using the training loop defined above.

Again, you don't need to tune any hyperparameters here, but you shoud achieve above 40% accuracy after one epoch of training.

```
[19]: # We need to wrap `flatten` function in a module in order to stack it
# in nn.Sequential
class Flatten(nn.Module):
    def forward(self, x):
        return flatten(x)

hidden_layer_size = 4000
learning_rate = 1e-2

model = nn.Sequential(
    Flatten(),
    nn.Linear(3 * 32 * 32, hidden_layer_size),
    nn.ReLU(),
    nn.Linear(hidden_layer_size, 10),
)

# you can use Nesterov momentum in optim.SGD
optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                      momentum=0.9, nesterov=True)

train_part34(model, optimizer)
```

```
Iteration 0, loss = 2.2854
Checking accuracy on validation set
Got 153 / 1000 correct (15.30)
```

```
Iteration 100, loss = 1.6150
Checking accuracy on validation set
Got 376 / 1000 correct (37.60)
```

```
Iteration 200, loss = 1.4701
Checking accuracy on validation set
Got 425 / 1000 correct (42.50)
```

```
Iteration 300, loss = 1.7515
Checking accuracy on validation set
Got 439 / 1000 correct (43.90)
```

```
Iteration 400, loss = 1.7775
Checking accuracy on validation set
Got 391 / 1000 correct (39.10)
```

```
Iteration 500, loss = 1.6248
Checking accuracy on validation set
Got 430 / 1000 correct (43.00)
```

```
Iteration 600, loss = 1.7612
Checking accuracy on validation set
Got 443 / 1000 correct (44.30)
```

```
Iteration 700, loss = 1.6102
Checking accuracy on validation set
Got 453 / 1000 correct (45.30)
```

## 7.0.2 Sequential API: Three-Layer ConvNet

Here you should use `nn.Sequential` to define and train a three-layer ConvNet with the same architecture we used in Part III:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You can use the default PyTorch weight initialization.

You should optimize your model using stochastic gradient descent with Nesterov momentum 0.9.

Again, you don't need to tune any hyperparameters but you should see accuracy above 55% after one epoch of training.

```
[20]: channel_1 = 32
channel_2 = 16
learning_rate = 7e-3

model = None
optimizer = None
```

```

#####
# TODO: Rewrite the 2-layer ConvNet with bias from Part III with the      #
# Sequential API.                                                       #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****#
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****#
#               END OF YOUR CODE          #
#####

train_part34(model, optimizer)

```

Iteration 0, loss = 2.3187  
 Checking accuracy on validation set  
 Got 123 / 1000 correct (12.30)

Iteration 100, loss = 1.3436  
 Checking accuracy on validation set  
 Got 496 / 1000 correct (49.60)

Iteration 200, loss = 1.1732  
 Checking accuracy on validation set  
 Got 488 / 1000 correct (48.80)

Iteration 300, loss = 1.5286  
 Checking accuracy on validation set  
 Got 524 / 1000 correct (52.40)

Iteration 400, loss = 1.1932  
 Checking accuracy on validation set  
 Got 532 / 1000 correct (53.20)

Iteration 500, loss = 1.4168  
 Checking accuracy on validation set  
 Got 498 / 1000 correct (49.80)

```
Iteration 600, loss = 1.1315
Checking accuracy on validation set
Got 554 / 1000 correct (55.40)
```

```
Iteration 700, loss = 1.0861
Checking accuracy on validation set
Got 563 / 1000 correct (56.30)
```

## 8 Part V. CIFAR-10 open-ended challenge

In this section, you can experiment with whatever ConvNet architecture you'd like on CIFAR-10.

Now it's your job to experiment with architectures, hyperparameters, loss functions, and optimizers to train a model that achieves **at least 70%** accuracy on the CIFAR-10 **validation** set within 10 epochs. You can use the `check_accuracy` and `train` functions from above. You can use either `nn.Module` or `nn.Sequential` API.

Describe what you did at the end of this notebook.

Here are the official API documentation for each component. One note: what we call in the class "spatial batch norm" is called "BatchNorm2D" in PyTorch.

- Layers in `torch.nn` package: <http://pytorch.org/docs/stable/nn.html>
- Activations: <http://pytorch.org/docs/stable/nn.html#non-linear-activations>
- Loss functions: <http://pytorch.org/docs/stable/nn.html#loss-functions>
- Optimizers: <http://pytorch.org/docs/stable/optim.html>

### 8.0.1 Things you might try:

- **Filter size:** Above we used 5x5; would smaller filters be more efficient?
- **Number of filters:** Above we used 32 filters. Do more or fewer do better?
- **Pooling vs Strided Convolution:** Do you use max pooling or just stride convolutions?
- **Batch normalization:** Try adding spatial batch normalization after convolution layers and vanilla batch normalization after affine layers. Do your networks train faster?
- **Network architecture:** The network above has two layers of trainable parameters. Can you do better with a deep network? Good architectures to try include:
  - [conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
  - [conv-relu-conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
  - [batchnorm-relu-conv]xN -> [affine]xM -> [softmax or SVM]
- **Global Average Pooling:** Instead of flattening and then having multiple affine layers, perform convolutions until your image gets small (7x7 or so) and then perform an average pooling operation to get to a 1x1 image picture (1, 1 , Filter#), which is then reshaped into a (Filter#) vector. This is used in [Google's Inception Network](#) (See Table 1 for their architecture).
- **Regularization:** Add l2 weight regularization, or perhaps use Dropout.

## 8.0.2 Tips for training

For each network architecture that you try, you should tune the learning rate and other hyperparameters. When doing this there are a couple important things to keep in mind:

- If the parameters are working well, you should see improvement within a few hundred iterations
- Remember the coarse-to-fine approach for hyperparameter tuning: start by testing a large range of hyperparameters for just a few training iterations to find the combinations of parameters that are working at all.
- Once you have found some sets of parameters that seem to work, search more finely around these parameters. You may need to train for more epochs.
- You should use the validation set for hyperparameter search, and save your test set for evaluating your architecture on the best parameters as selected by the validation set.

## 8.0.3 Going above and beyond

If you are feeling adventurous there are many other features you can implement to try and improve your performance. You are **not required** to implement any of these, but don't miss the fun if you have time!

- Alternative optimizers: you can try Adam, Adagrad, RMSprop, etc.
- Alternative activation functions such as leaky ReLU, parametric ReLU, ELU, or MaxOut.
- Model ensembles
- Data augmentation
- New Architectures
  - ResNets where the input from the previous layer is added to the output.
  - DenseNets where inputs into previous layers are concatenated together.
  - This blog has an [in-depth overview](#)

## 8.0.4 Have fun and happy training!

```
[27] : #####  
# TODO:  
# Experiment with any architectures, optimizers, and hyperparameters.  
# Achieve AT LEAST 70% accuracy on the *validation set* within 10 epochs.  
#  
# Note that you can use the check_accuracy function to evaluate on either  
# the test set or the validation set, by passing either loader_test or  
# loader_val as the second argument to check_accuracy. You should not touch  
# the test set until you have finished your architecture and hyperparameter  
# tuning, and only run the test set once at the end to report a final value.  
#####  
model = None  
optimizer = None  
  
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****  
import torch  
import torch.nn as nn
```

```

import torch.nn.functional as F

class CifarNet(nn.Module):
    def __init__(self, in_channels, num_classes, conv_channels, affine_units,
     ↵N=2, M=2, dropout_prob=0.6):
        super(CifarNet, self).__init__()

        self.conv_blocks = nn.ModuleList()
        for i in range(N):
            in_ch = in_channels if i == 0 else conv_channels[i - 1]
            out_ch = conv_channels[i]
            self.conv_blocks.append(nn.Sequential(
                nn.BatchNorm2d(in_ch),
                nn.ReLU(),
                nn.Conv2d(in_ch, out_ch, kernel_size=3, stride=1, padding=1)
            ))

        self.fc_layers = nn.ModuleList()
        for j in range(M):
            in_features = conv_channels[-1] * 32 * 32 if j == 0 else
     ↵affine_units[j - 1]
            out_features = affine_units[j]
            self.fc_layers.append(nn.Sequential(
                nn.Linear(in_features, out_features),
                nn.Dropout(dropout_prob),
                nn.ReLU()
            ))

        self.classifier = nn.Linear(affine_units[-1], num_classes)

    def forward(self, x):
        for conv_block in self.conv_blocks:
            x = conv_block(x)

        x = x.view(x.size(0), -1)

        for fc_layer in self.fc_layers:
            x = fc_layer(x)

        x = self.classifier(x)

    return x

model = CifarNet(
    in_channels=3,
    num_classes=10,
    conv_channels=[16, 32, 64, 32],

```

```

    affine_units=[512, 256],
    N=4,
    M=2,
)

optimizer = optim.Adam(model.parameters(), lr=0.001)

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

# You should get at least 70% accuracy
train_part34(model, optimizer, epochs=10)

```

Iteration 0, loss = 2.3221  
 Checking accuracy on validation set  
 Got 153 / 1000 correct (15.30)

Iteration 100, loss = 2.0317  
 Checking accuracy on validation set  
 Got 298 / 1000 correct (29.80)

Iteration 200, loss = 1.9956  
 Checking accuracy on validation set  
 Got 309 / 1000 correct (30.90)

Iteration 300, loss = 1.8907  
 Checking accuracy on validation set  
 Got 366 / 1000 correct (36.60)

Iteration 400, loss = 2.0118  
 Checking accuracy on validation set  
 Got 400 / 1000 correct (40.00)

Iteration 500, loss = 1.9244  
 Checking accuracy on validation set  
 Got 432 / 1000 correct (43.20)

Iteration 600, loss = 1.6181  
 Checking accuracy on validation set  
 Got 447 / 1000 correct (44.70)

Iteration 700, loss = 1.6102  
 Checking accuracy on validation set  
 Got 473 / 1000 correct (47.30)

Iteration 0, loss = 1.5243  
Checking accuracy on validation set  
Got 452 / 1000 correct (45.20)

Iteration 100, loss = 1.8411  
Checking accuracy on validation set  
Got 455 / 1000 correct (45.50)

Iteration 200, loss = 1.6200  
Checking accuracy on validation set  
Got 513 / 1000 correct (51.30)

Iteration 300, loss = 1.4919  
Checking accuracy on validation set  
Got 461 / 1000 correct (46.10)

Iteration 400, loss = 1.5268  
Checking accuracy on validation set  
Got 487 / 1000 correct (48.70)

Iteration 500, loss = 1.6984  
Checking accuracy on validation set  
Got 493 / 1000 correct (49.30)

Iteration 600, loss = 1.5203  
Checking accuracy on validation set  
Got 491 / 1000 correct (49.10)

Iteration 700, loss = 1.5870  
Checking accuracy on validation set  
Got 491 / 1000 correct (49.10)

Iteration 0, loss = 1.3685  
Checking accuracy on validation set  
Got 528 / 1000 correct (52.80)

Iteration 100, loss = 1.5175  
Checking accuracy on validation set  
Got 554 / 1000 correct (55.40)

Iteration 200, loss = 1.5691  
Checking accuracy on validation set  
Got 584 / 1000 correct (58.40)

Iteration 300, loss = 1.4291  
Checking accuracy on validation set  
Got 568 / 1000 correct (56.80)

Iteration 400, loss = 1.2851  
Checking accuracy on validation set  
Got 554 / 1000 correct (55.40)

Iteration 500, loss = 1.2842  
Checking accuracy on validation set  
Got 579 / 1000 correct (57.90)

Iteration 600, loss = 1.3411  
Checking accuracy on validation set  
Got 601 / 1000 correct (60.10)

Iteration 700, loss = 1.3074  
Checking accuracy on validation set  
Got 596 / 1000 correct (59.60)

Iteration 0, loss = 1.0189  
Checking accuracy on validation set  
Got 620 / 1000 correct (62.00)

Iteration 100, loss = 1.3212  
Checking accuracy on validation set  
Got 625 / 1000 correct (62.50)

Iteration 200, loss = 0.9773  
Checking accuracy on validation set  
Got 629 / 1000 correct (62.90)

Iteration 300, loss = 1.2878  
Checking accuracy on validation set  
Got 652 / 1000 correct (65.20)

Iteration 400, loss = 1.2292  
Checking accuracy on validation set  
Got 643 / 1000 correct (64.30)

Iteration 500, loss = 1.4515  
Checking accuracy on validation set  
Got 598 / 1000 correct (59.80)

Iteration 600, loss = 1.2565  
Checking accuracy on validation set  
Got 644 / 1000 correct (64.40)

Iteration 700, loss = 0.9177  
Checking accuracy on validation set  
Got 634 / 1000 correct (63.40)

Iteration 0, loss = 1.1416  
Checking accuracy on validation set  
Got 652 / 1000 correct (65.20)

Iteration 100, loss = 1.1849  
Checking accuracy on validation set  
Got 656 / 1000 correct (65.60)

Iteration 200, loss = 0.8104  
Checking accuracy on validation set  
Got 652 / 1000 correct (65.20)

Iteration 300, loss = 0.8795  
Checking accuracy on validation set  
Got 669 / 1000 correct (66.90)

Iteration 400, loss = 0.9805  
Checking accuracy on validation set  
Got 654 / 1000 correct (65.40)

Iteration 500, loss = 1.0973  
Checking accuracy on validation set  
Got 686 / 1000 correct (68.60)

Iteration 600, loss = 0.9939  
Checking accuracy on validation set  
Got 672 / 1000 correct (67.20)

Iteration 700, loss = 1.0197  
Checking accuracy on validation set  
Got 689 / 1000 correct (68.90)

Iteration 0, loss = 0.8916  
Checking accuracy on validation set  
Got 667 / 1000 correct (66.70)

Iteration 100, loss = 0.9392  
Checking accuracy on validation set  
Got 700 / 1000 correct (70.00)

Iteration 200, loss = 0.9835  
Checking accuracy on validation set  
Got 698 / 1000 correct (69.80)

Iteration 300, loss = 1.0694  
Checking accuracy on validation set  
Got 680 / 1000 correct (68.00)

Iteration 400, loss = 1.0011  
Checking accuracy on validation set  
Got 716 / 1000 correct (71.60)

Iteration 500, loss = 0.8417  
Checking accuracy on validation set  
Got 686 / 1000 correct (68.60)

Iteration 600, loss = 1.0643  
Checking accuracy on validation set  
Got 718 / 1000 correct (71.80)

Iteration 700, loss = 0.6745  
Checking accuracy on validation set  
Got 702 / 1000 correct (70.20)

Iteration 0, loss = 0.9069  
Checking accuracy on validation set  
Got 698 / 1000 correct (69.80)

Iteration 100, loss = 0.7527  
Checking accuracy on validation set  
Got 716 / 1000 correct (71.60)

Iteration 200, loss = 0.9449  
Checking accuracy on validation set  
Got 697 / 1000 correct (69.70)

Iteration 300, loss = 0.9958  
Checking accuracy on validation set  
Got 702 / 1000 correct (70.20)

Iteration 400, loss = 0.9136  
Checking accuracy on validation set  
Got 669 / 1000 correct (66.90)

Iteration 500, loss = 1.0299  
Checking accuracy on validation set  
Got 713 / 1000 correct (71.30)

Iteration 600, loss = 0.8024  
Checking accuracy on validation set  
Got 742 / 1000 correct (74.20)

Iteration 700, loss = 0.7521  
Checking accuracy on validation set  
Got 711 / 1000 correct (71.10)

```
Iteration 0, loss = 0.8735
Checking accuracy on validation set
Got 720 / 1000 correct (72.00)
```

```
Iteration 100, loss = 0.7720
Checking accuracy on validation set
Got 713 / 1000 correct (71.30)
```

```
Iteration 200, loss = 0.5504
Checking accuracy on validation set
Got 711 / 1000 correct (71.10)
```

```
Iteration 300, loss = 0.8049
Checking accuracy on validation set
Got 720 / 1000 correct (72.00)
```

```
Iteration 400, loss = 0.7905
Checking accuracy on validation set
Got 709 / 1000 correct (70.90)
```

```
Iteration 500, loss = 0.7352
Checking accuracy on validation set
Got 719 / 1000 correct (71.90)
```

```
Iteration 600, loss = 0.9139
Checking accuracy on validation set
Got 715 / 1000 correct (71.50)
```

```
Iteration 700, loss = 0.8139
Checking accuracy on validation set
Got 712 / 1000 correct (71.20)
```

```
Iteration 0, loss = 0.6837
Checking accuracy on validation set
Got 712 / 1000 correct (71.20)
```

```
Iteration 100, loss = 0.6726
Checking accuracy on validation set
Got 723 / 1000 correct (72.30)
```

```
Iteration 200, loss = 0.4234
Checking accuracy on validation set
Got 733 / 1000 correct (73.30)
```

```
Iteration 300, loss = 0.7471
Checking accuracy on validation set
Got 703 / 1000 correct (70.30)
```

```
Iteration 400, loss = 1.0109
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)
```

```
Iteration 500, loss = 1.0142
Checking accuracy on validation set
Got 733 / 1000 correct (73.30)
```

```
Iteration 600, loss = 0.8164
Checking accuracy on validation set
Got 739 / 1000 correct (73.90)
```

```
Iteration 700, loss = 0.4185
Checking accuracy on validation set
Got 741 / 1000 correct (74.10)
```

```
Iteration 0, loss = 0.5374
Checking accuracy on validation set
Got 724 / 1000 correct (72.40)
```

```
Iteration 100, loss = 0.5396
Checking accuracy on validation set
Got 731 / 1000 correct (73.10)
```

```
Iteration 200, loss = 0.6241
Checking accuracy on validation set
Got 761 / 1000 correct (76.10)
```

```
Iteration 300, loss = 0.5267
Checking accuracy on validation set
Got 750 / 1000 correct (75.00)
```

```
Iteration 400, loss = 0.8058
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)
```

```
Iteration 500, loss = 0.5851
Checking accuracy on validation set
Got 735 / 1000 correct (73.50)
```

```
Iteration 600, loss = 0.5350
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)
```

```
Iteration 700, loss = 0.6672
Checking accuracy on validation set
Got 729 / 1000 correct (72.90)
```

## 8.1 Describe what you did

In the cell below you should write an explanation of what you did, any additional features that you implemented, and/or any graphs that you made in the process of training and evaluating your network.

**Answer:**

## 8.2 Test set – run this only once

Now that we've gotten a result we're happy with, we test our final model on the test set (which you should store in best\_model). Think about how this compares to your validation set accuracy.

```
[28]: best_model = model  
check_accuracy_part34(loader_test, best_model)
```

```
Checking accuracy on test set  
Got 7298 / 10000 correct (72.98)
```

```
[22]:
```

# Network\_Visualization

November 14, 2024

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cse493g1/assignments/assignment4/'
FOLDERNAME = 'cse493g1/assignments/assignment4'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My\ Drive/{}'.format(FOLDERNAME))

# This downloads the COCO dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cse493g1/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My\ Drive/cse493g1/assignments/assignment4/cse493g1/datasets
/content/drive/My\ Drive/cse493g1/assignments/assignment4
```

## 1 Network Visualization

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 CNN 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 we will 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.

We will explore three techniques for image generation.

**Saliency Maps.** We can use saliency maps to tell which part of the image influenced the classification decision made by the network.

**Fooling Images.** We can perturb an input image so that it appears the same to humans but will be misclassified by the pretrained network.

**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.

```
[2]: # Setup cell.  
import torch  
import torchvision  
import numpy as np  
import random  
import matplotlib.pyplot as plt  
from PIL import Image  
from cse493g1.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD  
from cse493g1.net_visualization_pytorch import *  
  
%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'  
  
%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.

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

```
[3]: # Download and load the pretrained SqueezeNet model.  
model = torchvision.models.squeezenet1_1(pretrained=True)  
  
# We don't want to train the model, so tell PyTorch not to compute gradients
```

```
# with respect to model parameters.
for param in model.parameters():
    param.requires_grad = False
```

```
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:208:
UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be
removed in the future, please use 'weights' instead.
    warnings.warn(
/usr/local/lib/python3.10/dist-packages/torchvision/models/_utils.py:223:
UserWarning: Arguments other than a weight enum or `None` for 'weights' are
deprecated since 0.13 and may be removed in the future. The current behavior is
equivalent to passing `weights=SqueezeNet1_1_Weights.IMGNET1K_V1`. You can
also use `weights=SqueezeNet1_1_Weights.DEFAULT` to get the most up-to-date
weights.
    warnings.warn(msg)
Downloading: "https://download.pytorch.org/models/squeezeimagenet1_1-b8a52dc0.pth" to
/root/.cache/torch/hub/checkpoints/squeezeimagenet1_1-b8a52dc0.pth
100%| 4.73M/4.73M [00:00<00:00, 18.1MB/s]
```

## 2.1 Loading ImageNet Validation Images

We have provided a few example images from the validation set of the ImageNet ILSVRC 2012 Classification dataset. 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.

```
[4]: from cse493g1.data_utils import load_imagenet_val
X, 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[i])
    plt.title(class_names[y[i]])
    plt.axis('off')
plt.gcf().tight_layout()
```



### 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 (3, H, W) then this gradient will also have shape (3, H, W); for each pixel in the image, this gradient tells us the amount by which the classification score will change if the pixel changes by a small amount. To compute the saliency map, we take the absolute value of this gradient, then take the maximum value over the 3 input channels; the final saliency map thus has shape (H, W) and all entries are nonnegative.

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

#### 3.0.1 Hint: PyTorch gather 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  $s[\text{np.arange}(N), y]$  is a numpy array of shape (N,) which selects one element from each element in  $s$  using the indices in  $y$ .

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

```
s.gather(1, y.view(-1, 1)).squeeze()
```

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

run the following cell to see an example.

You can also read the documentation for [the gather method](#) and [the squeeze method](#).

[5]: # Example of using gather to select one entry from each row in PyTorch

```
def gather_example():
    N, C = 4, 5
    s = torch.randn(N, C)
    y = torch.LongTensor([1, 2, 1, 3])
    print(s)
    print(y)
    print(s.gather(1, y.view(-1, 1)).squeeze())
gather_example()
```

```
tensor([[ 2.2521,   0.2705,  -0.5338,   1.4999,   2.2990],
        [ 0.1118,  -0.3538,  -0.7881,   0.4363,  -0.5219],
        [-0.4955,   1.1380,  -1.8610,   0.0453,   0.6132],
        [-1.1333,  -0.8946,  -0.8791,  -0.2775,   0.8924]])
tensor([1, 2, 1, 3])
tensor([ 0.2705, -0.7881,  1.1380, -0.2775])
```

Implement `compute_saliency_maps` function inside `cse493g1/net_visualization_pytorch.py`

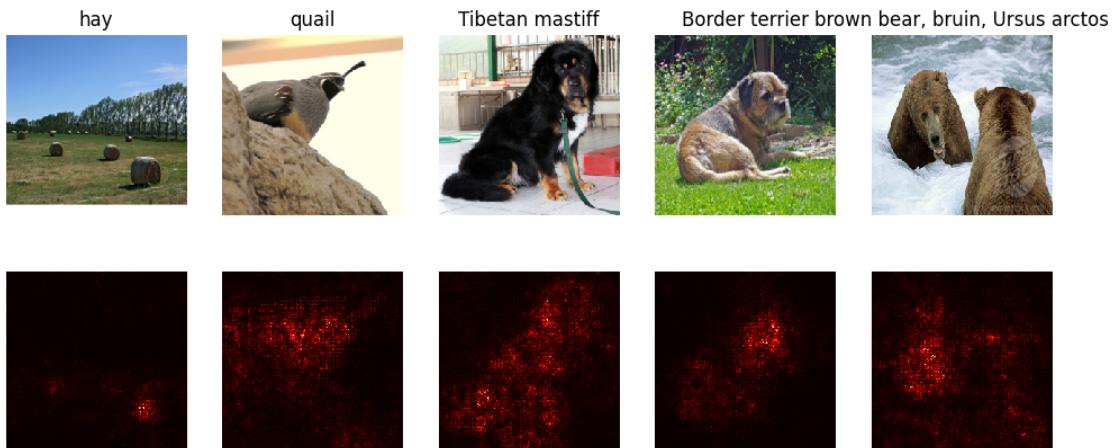
Once you have completed the implementation above, run the following to visualize some class saliency maps on our example images from the ImageNet validation set:

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

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

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

show_saliency_maps(X, y)
```



## 4 Inline Question 1

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

**Your Answer:** No, this assertion is false because : 1. Saliency maps store the maximum absolute value of the gradient across the color channels, but due to this, they are not directional. Therefore, we would likely end up with random or oscillating updates, as each pixel would change without regard to the correct direction. 2. Further, when we take the maximum gradient across channels (in the case of an RGB image), we club information for all the colours together and lose information about which specific color channel is more important for the score. Gradient ascent requires this detailed information to update each channel individually in the direction that will increase the class score.

## 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 `cse493g1/net_visualization_pytorch.py`

Run the following cell 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.

```
[7]: idx = 0
target_y = 6

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

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

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

```
[8]: X_fooling_np = deprocess(X_fooling.clone())
X_fooling_np = np.asarray(X_fooling_np).astype(np.uint8)

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

```

plt.axis('off')

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

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

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

plt.gcf().set_size_inches(12, 5)
plt.show()

```



## 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

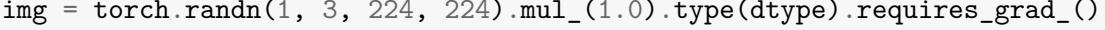
In `cse493g1/net_visualization_pytorch.py` complete the implementation of the `class_visualization_update_step` used in the `create_class_visualization` function below. Once you have completed that implementation, run the following cells to generate an image of a Tarantula:

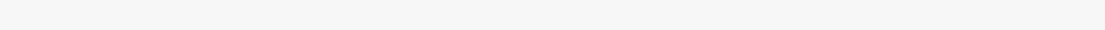
```
[14]: def create_class_visualization(target_y, model, dtype, **kwargs):
    """
    Generate an image to maximize the score of target_y under a pretrained
    ↵model.

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

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

```

# Randomly initialize the image as a PyTorch Tensor, and make it requires_grad_.
 gradient.

 = torch.randn(1, 3, 224, 224).mul_(1.0).type(dtype).requires_grad_()

for t in range(num_iterations):
    # Randomly jitter the image a bit; this gives slightly nicer results
    ox, oy = random.randint(0, max_jitter), random.randint(0, max_jitter)
    img.data.copy_(jitter(img.data, ox, oy))
    class_visualization_update_step(img, model, target_y, l2_reg, learning_rate)
    # Undo the random jitter
    img.data.copy_(jitter(img.data, -ox, -oy))

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

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

    return deprocess(img.data.cpu())

```

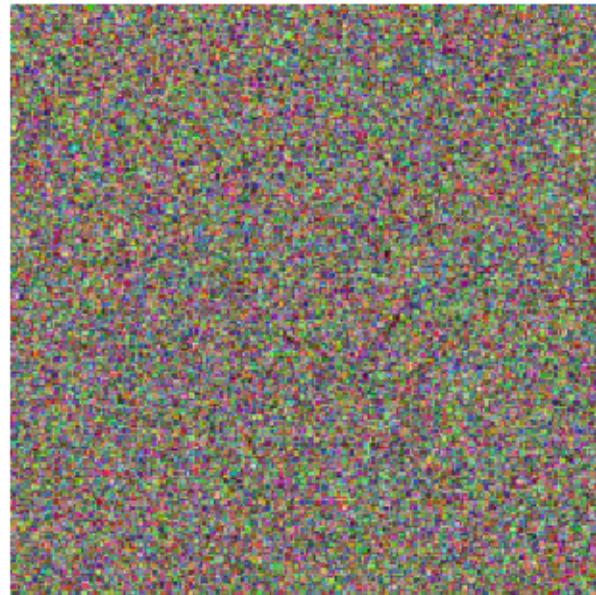
```

[15]: dtype = torch.FloatTensor
model.type(dtype)

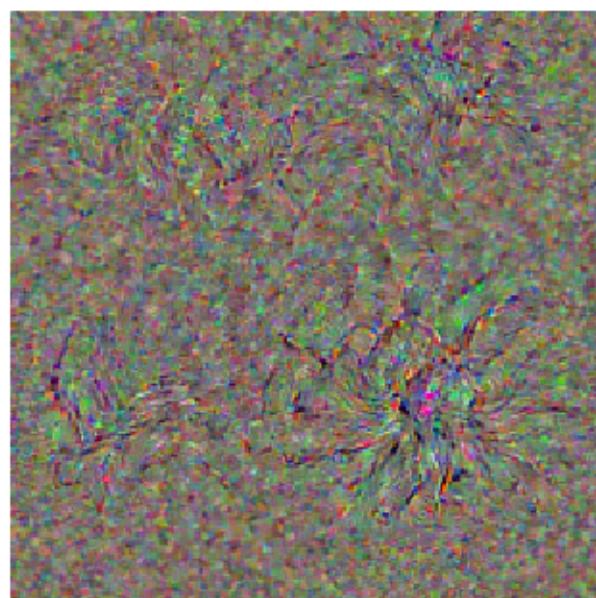
target_y = 76 # Tarantula
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
out = create_class_visualization(target_y, model, dtype)

```

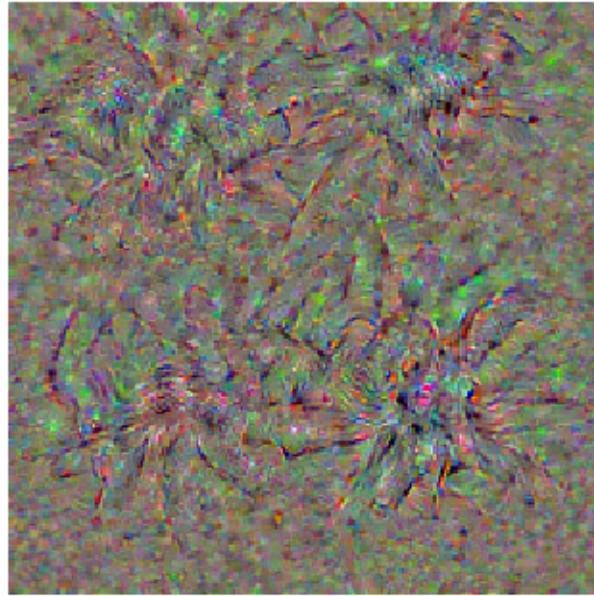
tarantula  
Iteration 1 / 100



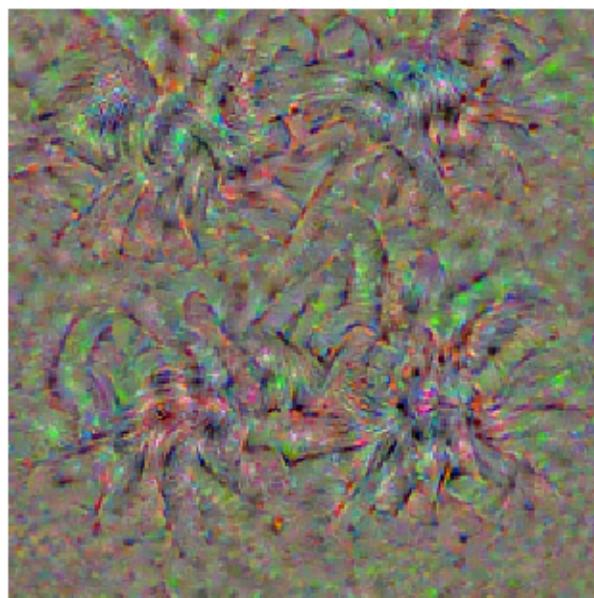
tarantula  
Iteration 25 / 100

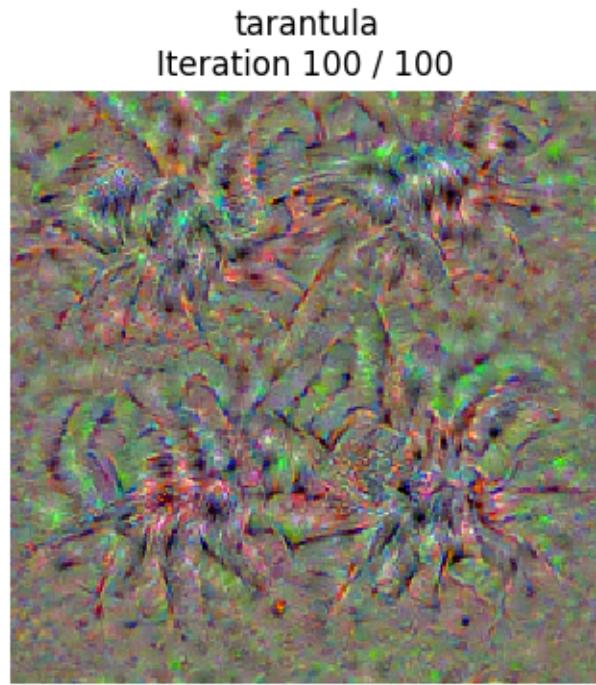


tarantula  
Iteration 50 / 100



tarantula  
Iteration 75 / 100



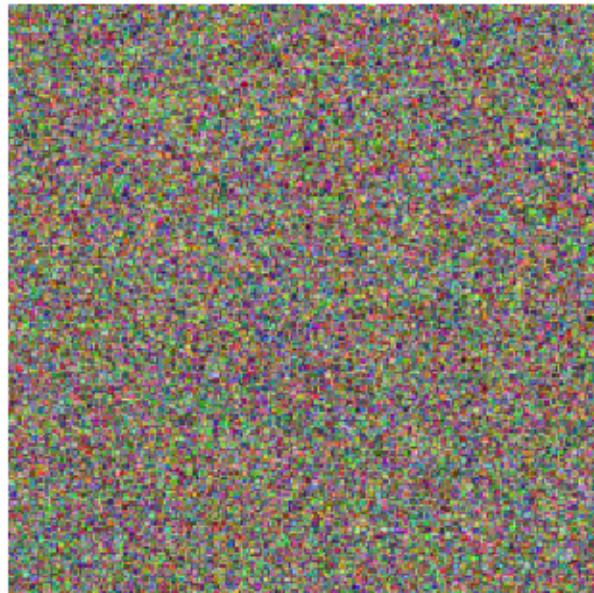


Try out your class visualization on other classes! You should also feel free to play with various hyperparameters to try and improve the quality of the generated image, but this is not required.

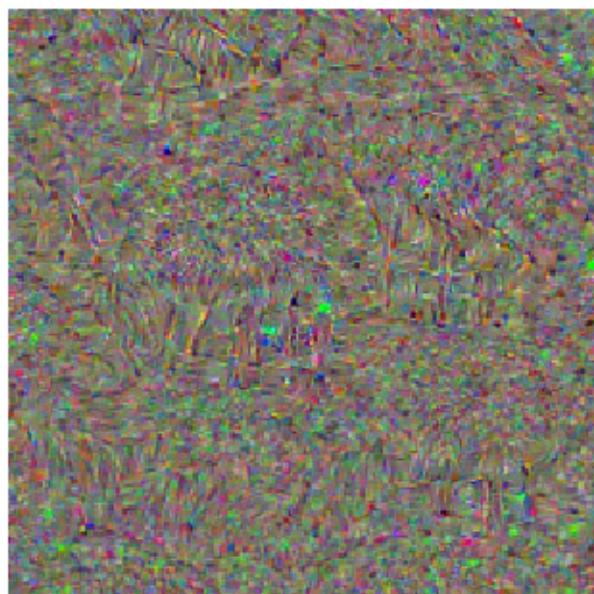
```
[16]: # target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
target_y = np.random.randint(1000)
print(class_names[target_y])
X = create_class_visualization(target_y, model, dtype)
```

thatch, thatched roof

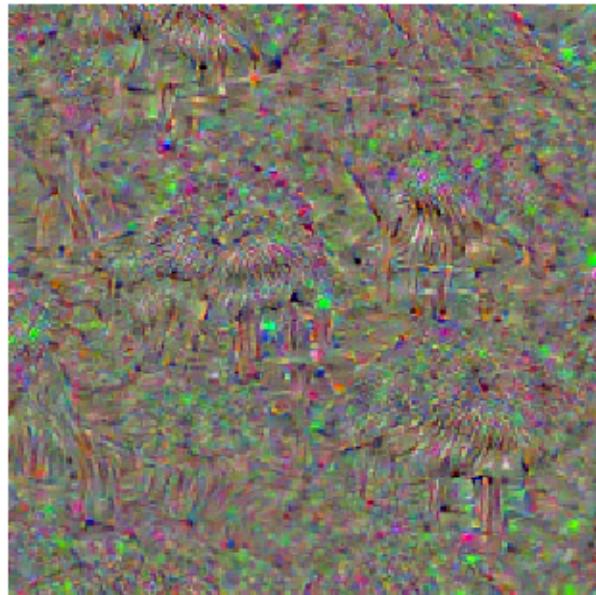
thatch, thatched roof  
Iteration 1 / 100



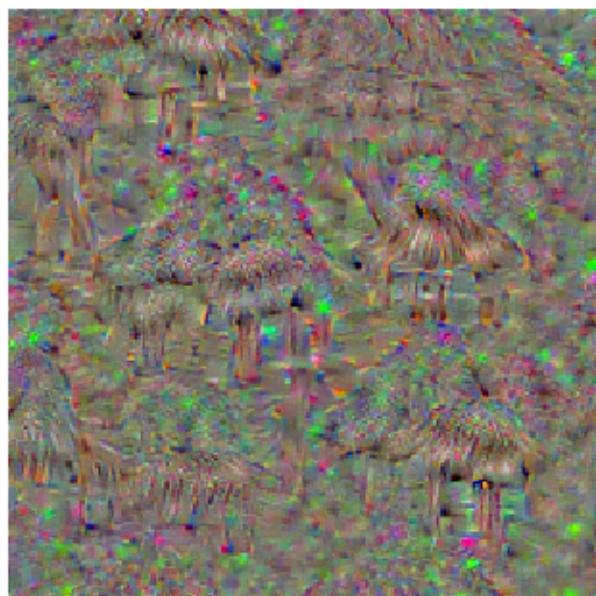
thatch, thatched roof  
Iteration 25 / 100



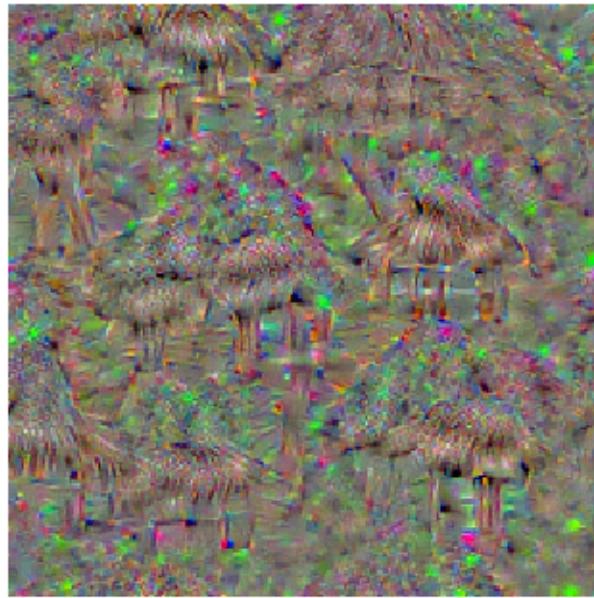
thatch, thatched roof  
Iteration 50 / 100



thatch, thatched roof  
Iteration 75 / 100



thatch, thatched roof  
Iteration 100 / 100



# RNN\_Captioning

November 14, 2024

```
[1]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cse493g1/assignments/assignment4/'
FOLDERNAME = 'cse493g1/assignments/assignment4'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My\ Drive/{}'.format(FOLDERNAME))

# This downloads the COCO dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cse493g1/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

Mounted at /content/drive  
/content/drive/My Drive/cse493g1/assignments/assignment4/cse493g1/datasets  
/content/drive/My Drive/cse493g1/assignments/assignment4

## 1 Image Captioning with RNNs

In this exercise, you will implement vanilla Recurrent Neural Networks and use them to train a model that can generate novel captions for images.

```
[2]: # Setup cell.
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cse493g1.gradient_check import eval_numerical_gradient, ▾
    eval_numerical_gradient_array
```

```

from cse493g1.rnn_layers import *
from cse493g1.captioning_solver import CaptioningSolver
from cse493g1.classifiers.rnn import CaptioningRNN
from cse493g1.coco_utils import load_coco_data, sample_coco_minibatch, decode_captions
from cse493g1.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'

%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 COCO Dataset

For this exercise, we will use the 2014 release of the [COCO dataset](#), a 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.

**Image features.** 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, and these features are stored in the files `train2014_vgg16_fc7.h5` and `val2014_vgg16_fc7.h5`. To cut down on processing time and memory requirements, we have reduced the dimensionality of the features from 4096 to 512 using Principal Component Analysis (PCA), and these features are stored in the files `train2014_vgg16_fc7_pca.h5` and `val2014_vgg16_fc7_pca.h5`. The raw images take up nearly 20GB of space so we have not included them in the download. Since all images are taken from Flickr, we have stored the URLs of the training and validation images in the files `train2014_urls.txt` and `val2014_urls.txt`. This allows you to download images on-the-fly for visualization.

**Captions.** 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 `cse493g1/coco_utils.py` to convert NumPy arrays of integer IDs back into strings.

**Tokens.** There are a couple special tokens that we add to the vocabulary, and we have taken care of all implementation details around special tokens for you. 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.

You can load all of the COCO data (captions, features, URLs, and vocabulary) using the `load_coco_data` function from the file `cse493g1/coco_utils.py`. Run the following cell to do so:

```
[3]: # Load COCO data from disk into a dictionary.  
# We'll work with dimensionality-reduced features for the remainder of this ↴  
# assignment,  
# but you can also experiment with the original features on your own 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 /content/drive/My  
Drive/cse493g1/assignments/assignment4/cse493g1/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
```

## 2.1 Inspect 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 `cse493g1/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.

```
[4]: # Sample a minibatch and show the images and captions.  
# If you get an error, the URL just no longer exists, so don't worry!  
# You can re-sample as many times as you want.  
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 = decodeCaptions(caption, data['idx_to_word'])
plt.title(caption_str)
plt.show()

```

Output hidden; open in <https://colab.research.google.com> to view.

### 3 Recurrent Neural Network

As discussed in lecture, we will use Recurrent Neural Network (RNN) language models for image captioning. The file `cse493g1/rnn_layers.py` contains implementations of different layer types that are needed for recurrent neural networks, and the file `cse493g1/classifiers/rnn.py` uses these layers to implement an image captioning model.

We will first implement different types of RNN layers in `cse493g1/rnn_layers.py`.

**NOTE:** The Long-Short Term Memory (LSTM) RNN is a common variant of the vanilla RNN. `LSTM_Captioning.ipynb` is optional extra credit, so don't worry about references to LSTM in `cse493g1/classifiers/rnn.py` and `cse493g1/rnn_layers.py` for now.

### 4 Copy Optimizers from Previous Assignments

Copy the optimizers you have coded already into `cse493g1/optim.py`

### 5 Vanilla RNN: Step Forward

Open the file `cse493g1/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.

[5]: N, D, H = 3, 10, 4

```

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

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

```

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

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

## 6 Vanilla RNN: Step Backward

In the file `cse493g1/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.

```
[6]: from cse493g1.rnn_layers import rnn_step_forward, rnn_step_backward
np.random.seed(493)
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: 7.086878851139418e-10
dprev_h error: 1.6434946364267517e-10
dWx error: 1.8473161214202673e-09
dWh error: 2.9699750128300744e-09
db error: 2.2179207136496757e-10
```

## 7 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 `cse493g1/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.

```
[7]: 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
```

## 8 Vanilla RNN: Backward

In the file `cse493g1/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.

```
[8]: np.random.seed(493)
```

```
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:  1.8514653997303406e-09
dh0 error:  1.3756891368431857e-10
dWx error:  6.054940160837621e-09
dWh error:  2.4444390463236517e-10
db error:  1.3043315175619792e-09

```

## 9 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 `cse493g1/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.

```

[9]: N, T, V, D = 2, 4, 5, 3

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

out, _ = word_embedding_forward(x, W)
expected_out = np.asarray([
    [[ 0.,           0.07142857,  0.14285714], ...

```

```

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

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

```

out error: 1.000000094736443e-08

## 10 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.

```
[10]: np.random.seed(493)

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: 3.276698607443385e-12

## 11 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 `cse493g1/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.

```
[11]: np.random.seed(493)

# 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: 4.95166576746952e-11
dw error: 5.848067121783654e-11
db error: 3.902641368396469e-11

```

## 12 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 `cse493g1/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.

```
[12]: # Sanity check for temporal softmax loss
from cse493g1.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 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))

```

```

2.302616151141718
23.02577484954567
2.355991216891321
dx error: 5.844343451289026e-08

```

## 13 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 `cse493g1/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 on the order of  $e-10$  or less.

```
[13]: 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.

```

[14]: np.random.seed(493)

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: 5.541330e-08
W_proj relative error: 6.181084e-09
W_vocab relative error: 1.583895e-08
Wh relative error: 8.832173e-08
Wx relative error: 5.737853e-07
b relative error: 8.727016e-10
b_proj relative error: 1.410171e-09
b_vocab relative error: 1.536398e-10

```

## 14 Overfit RNN Captioning Model on 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 `cse493g1/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.

```
[15]: np.random.seed(493)

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()

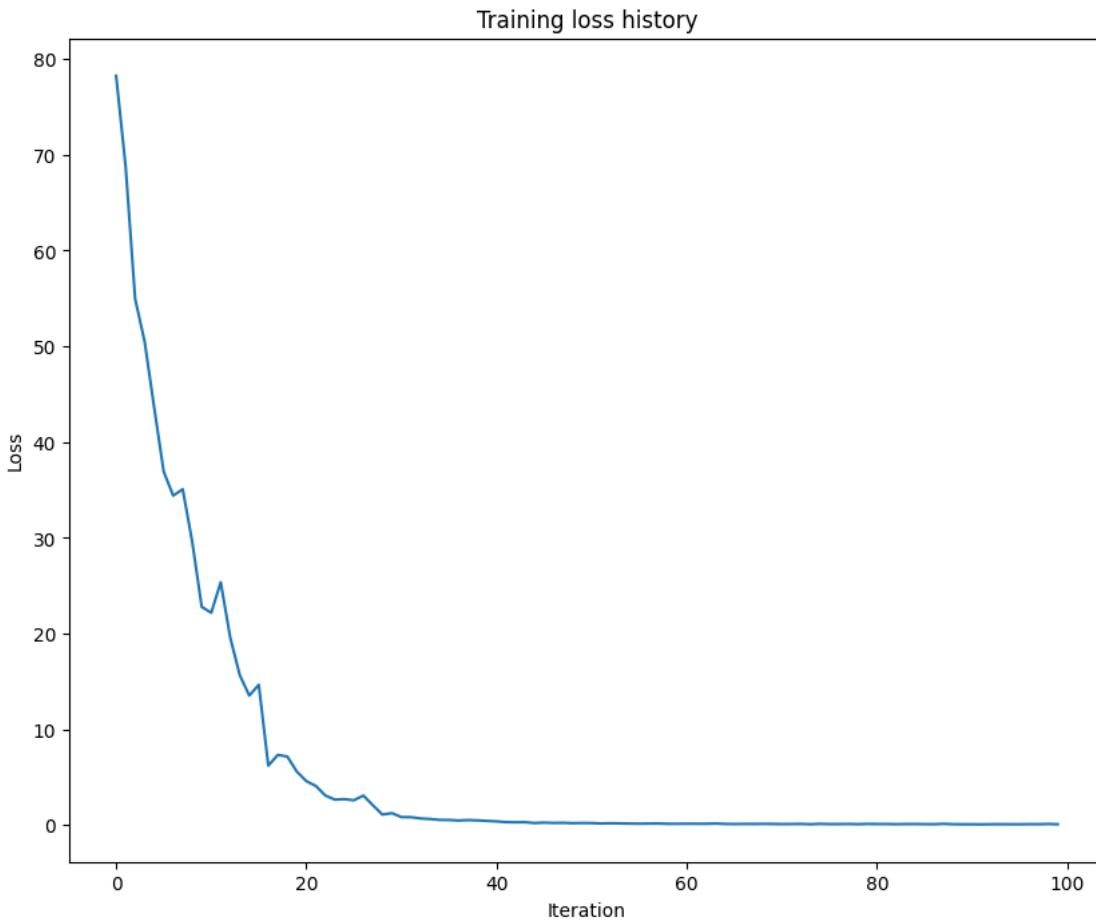
# Plot the training losses.
plt.plot(small_rnn_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()

```

```

base dir /content/drive/My
Drive/cse493g1/assignments/assignment4/cse493g1/datasets/coco_captioning
(Iteration 1 / 100) loss: 78.253366
(Iteration 11 / 100) loss: 22.183208
(Iteration 21 / 100) loss: 4.601543
(Iteration 31 / 100) loss: 0.842449
(Iteration 41 / 100) loss: 0.398038
(Iteration 51 / 100) loss: 0.220543
(Iteration 61 / 100) loss: 0.152443
(Iteration 71 / 100) loss: 0.114841
(Iteration 81 / 100) loss: 0.123294
(Iteration 91 / 100) loss: 0.093415

```



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

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

```
Final loss:  0.09787612782864435
```

## 15 RNN Sampling at Test Time

Unlike classification models, image captioning models behave very differently at training time vs. 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 `cse493g1/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, however, probably won't make sense.

```
[17]: # If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
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):
        img = image_from_url(url)
        # Skip missing URLs.
        if img is None: continue
        plt.imshow(img)
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
        plt.axis('off')
        plt.show()
```

train

a group of people standing around a living room near a couch <END>  
GT:<START> a group of people standing around a living room near a couch <END>



train

a yellow truck driving down a city street <END>  
GT:<START> a yellow truck driving down a city street <END>



val

several yellow eat piece <UNK> food in a mouth <END>  
GT:<START> two green trains stop in front of a <UNK> <UNK> <END>



val  
a man stands sitting on front of red <UNK> game <END>  
GT:<START> a man is watching something on a big screen <END>



## 16 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:** 1 clear advantage of an image captioning model that uses a character-level RNN is that it would have a far smaller vocabulary (and therefore, smaller embedding size) than a word-level RNN and it can be used to generate any word and would therefore, not suffer from the “UNKWN” token.

1 clear disadvantage, however, is that the sequence length of the caption would be far higher in the

character-level model than the word-level one, and so, it would take much longer (that is, many more timesteps) to generate the caption than the word-level model.

[17] :

# LSTM\_Captioning

November 14, 2024

```
[24]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cse493g1/assignments/assignment4/'
FOLDERNAME = 'cse493g1/assignments/assignment4'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My\ Drive/{}'.format(FOLDERNAME))

# This downloads the COCO dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cse493g1/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call  
drive.mount("/content/drive", force\_remount=True).  
/content/drive/My\ Drive/cse493g1/assignments/assignment4/cse493g1/datasets  
/content/drive/My\ Drive/cse493g1/assignments/assignment4

## 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.

```
[25]: # Setup cell.
import time, os, json
import numpy as np
import matplotlib.pyplot as plt
```

```

from cse493g1.gradient_check import eval_numerical_gradient,
    eval_numerical_gradient_array
from cse493g1.rnn_layers import *
from cse493g1.captioning_solver import CaptioningSolver
from cse493g1.classifiers.rnn import CaptioningRNN
from cse493g1.coco_utils import load_coco_data, sample_coco_minibatch,
    decode_captions
from cse493g1.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'

%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
```

## 2 COCO Dataset

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

```
[26]: # Load COCO data from disk into a dictionary.
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  /content/drive/MyDrive/cse493g1/assignments/assignment4/cse493g1/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
idxToWord <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

A common variant on the vanilla RNN is the Long-Short Term Memory (LSTM) RNN. 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  $\sigma$  is the sigmoid function and  $\tanh$  is the hyperbolic tangent, both applied elementwise.

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

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

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 `cse493g1/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.

```
[27]: N, D, H = 3, 4, 5
x = np.linspace(-0.4, 1.2, num=N*D).reshape(N, D)
prev_h = np.linspace(-0.3, 0.7, num=N*H).reshape(N, H)
```

```

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

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

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

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

```

```

next_h error:  5.7054131967097955e-09
next_c error:  5.8143123088804145e-09

```

## 5 LSTM: Step Backward

Implement the backward pass for a single LSTM timestep in the function `lstm_step_backward` in the file `cse493g1/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.

```

[28]: np.random.seed(493)

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: 4.3913766373906044e-10
dh error: 1.5384126336638421e-10
dc error: 1.1296260850605455e-09
dWx error: 5.133706248417227e-09
dWh error: 1.2463932282997943e-08
db error: 6.02569709548372e-09

```

## 6 LSTM: Forward

In the function `lstm_forward` in the file `cse493g1/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.

```

[29]: 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: 8.610537442272635e-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 `cse493g1/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).

```
[30]: from cse493g1.rnn_layers import lstm_forward, lstm_backward
np.random.seed(493)

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: 6.97770005072807e-09
dh0 error: 3.625318623162889e-09
dWx error: 2.767916213255284e-08
dWh error: 5.6706834117042436e-08
db error: 8.20572566716692e-10

```

## 8 LSTM Captioning Model

Now that you have implemented an LSTM, update the implementation of the `loss` method of the `CaptioningRNN` class in the file `cse493g1/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.

```

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

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

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

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

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

```

```

expected_loss = 9.82445935443

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

```

```

loss: 9.82445935443226
expected loss: 9.82445935443
difference: 2.261302256556519e-12

```

## 9 Overfit LSTM Captioning Model on Small Data

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.

```

[32]: np.random.seed(493)

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,
)

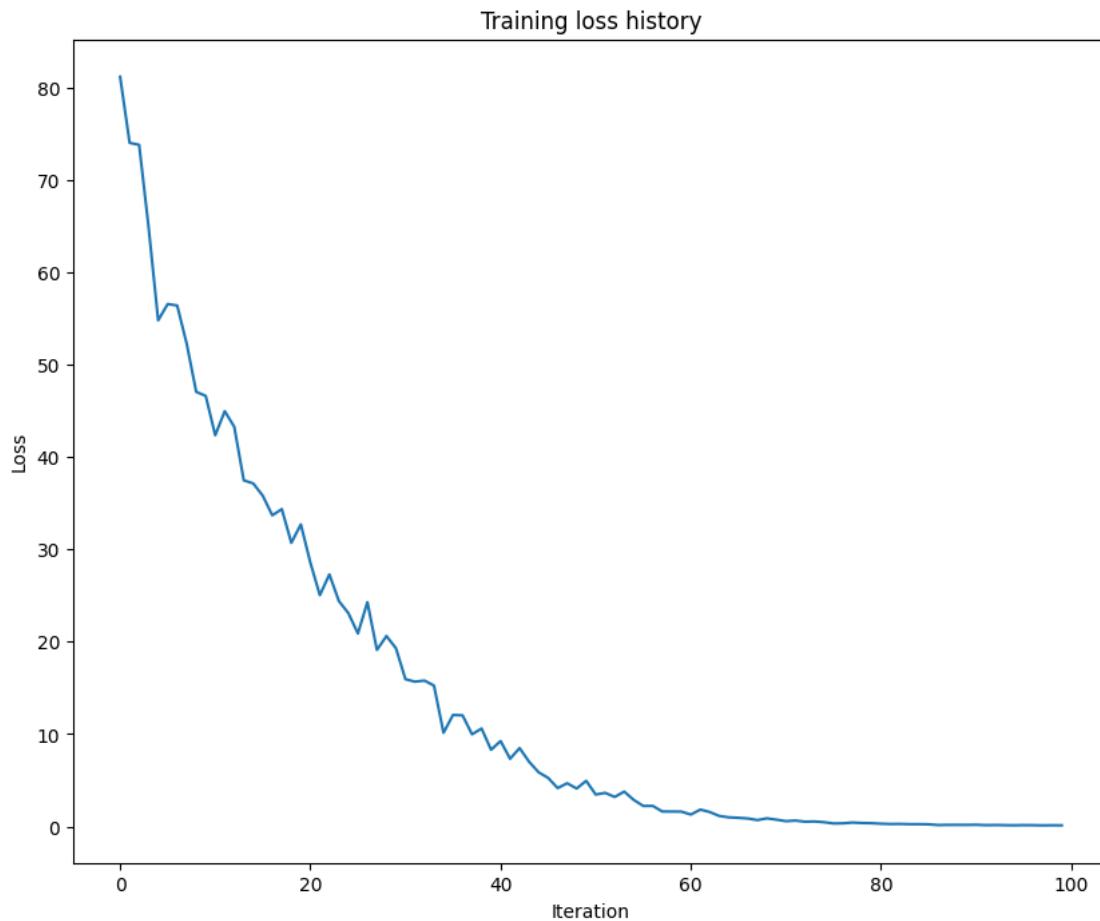
small_lstm_solver = CaptioningSolver(
    small_lstm_model, small_data,
    update_rule='adam',
    num_epochs=50,
    batch_size=25,
    optim_config={
        'learning_rate': 5e-3,
    },
    lr_decay=0.995,
    verbose=True, print_every=10,
)
small_lstm_solver.train()

# 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  /content/drive/MyDrive/cse493g1/assignments/assignment4/cse493g1/datasets/coco_captioning
(Iteration 1 / 100) loss: 81.125426
(Iteration 11 / 100) loss: 42.344705
(Iteration 21 / 100) loss: 28.563743
(Iteration 31 / 100) loss: 15.941721
(Iteration 41 / 100) loss: 9.255219
(Iteration 51 / 100) loss: 3.465271
(Iteration 61 / 100) loss: 1.298801
(Iteration 71 / 100) loss: 0.583298
(Iteration 81 / 100) loss: 0.307668
(Iteration 91 / 100) loss: 0.190979
```



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

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

```
Final loss: 0.13043152589277068
```

## 10 LSTM Sampling at Test Time

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).

```
[34]: # If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
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):
        img = image_from_url(url)
        # Skip missing URLs.
        if img is None: continue
        plt.imshow(img)
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
        plt.axis('off')
        plt.show()
```

train

this is a train going through a <UNK> area <END>  
GT:<START> this is a train going through a <UNK> area <END>



URL Error: Not Found

[http://farm9.staticflickr.com/8064/8171939431\\_07e742b75b\\_z.jpg](http://farm9.staticflickr.com/8064/8171939431_07e742b75b_z.jpg)

val

four guys playing a <UNK> bench on a grass <END>  
GT:<START> a big red tower on a landing by the ocean <END>



val

the young boy in a <UNK> is holding a baseball bat <END>  
GT:<START> a man holding a tennis racquet in his left hand <END>



[34] :