# Fully-Connected Neural Nets

In the previous homework you implemented a fully-connected two-layer neural network on CIFAR-10. The implementation was simple but not very modular since the loss and gradient were computed in a single monolithic function. This is manageable for a simple two-layer network, but would become impractical as we move to bigger models. Ideally we want to build networks using a more modular design so that we can implement different layer types in isolation and then snap them together into models with different architectures.

In this exercise we will implement fully-connected networks using a more modular approach. For each layer we will implement a `forward` and a `backward` function. The `forward` function will receive inputs, weights, and other parameters and will return both an output and a `cache` object storing data needed for the backward pass, like this:

```python
def layer_forward(x, w):
  """ Receive inputs x and weights w """
  # Do some computations ...
  z = # ... some intermediate value
  # Do some more computations ...
  out = # the output

  cache = (x, w, z, out) # Values we need to compute gradients

  return out, cache
```

The backward pass will receive upstream derivatives and the `cache` object, and will return gradients with respect to the inputs and weights, like this:

```python
def layer_backward(dout, cache):
  """
  Receive dout (derivative of loss with respect to outputs) and cache,
  and compute derivative with respect to inputs.
  """
  # Unpack cache values
  x, w, z, out = cache

  # Use values in cache to compute derivatives
  dx = # Derivative of loss with respect to x
  dw = # Derivative of loss with respect to w

  return dx, dw
```

After implementing a bunch of layers this way, we will be able to easily combine them to build classifiers with different architectures.

In addition to implementing fully-connected networks of arbitrary depth, we will also explore different update rules for optimization, and introduce Batch/Layer Normalization as a tool to more efficiently optimize deep networks.

In [ ]:
```python
# 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. 'cs6353/assignments/assignment3/'
FOLDERNAME = 'CS6353/Assignments/assignment3/assignment3/'
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 CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cs6353/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME

# Install requirements from colab_requirements.txt
# TODO: Please change your path below to the colab_requirements.txt file
! python -m pip install -r /content/drive/My\ Drive/$FOLDERNAME/requirements.txt
```

```
Mounted at /content/drive
/content/drive/My Drive/CS6353/Assignments/assignment3/assignment3/cs6353/datasets
--2024-10-18 16:18:31--  https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz
Resolving www.cs.toronto.edu (www.cs.toronto.edu)... 128.100.3.30
Connecting to www.cs.toronto.edu (www.cs.toronto.edu)|128.100.3.30|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 170498071 (163M) [application/x-gzip]
Saving to: 'cifar-10-python.tar.gz'

cifar-10-python.tar 100%[===================>] 162.60M  36.4MB/s    in 4.7s

2024-10-18 16:18:35 (34.5 MB/s) - 'cifar-10-python.tar.gz' saved [170498071/17049807
1]

cifar-10-batches-py/
cifar-10-batches-py/data_batch_4
cifar-10-batches-py/readme.html
cifar-10-batches-py/test_batch
cifar-10-batches-py/data_batch_3
cifar-10-batches-py/batches.meta
cifar-10-batches-py/data_batch_2
cifar-10-batches-py/data_batch_5
cifar-10-batches-py/data_batch_1
/content/drive/My Drive/CS6353/Assignments/assignment3/assignment3
Collecting attrs==19.1.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignme
nt3/assignment3//requirements.txt (line 1))
  Downloading attrs-19.1.0-py2.py3-none-any.whl.metadata (10 kB)
Collecting backcall==0.1.0 (from -r /content/drive/My Drive/CS6353/Assignments/assign
ment3/assignment3//requirements.txt (line 2))
  Downloading backcall-0.1.0.zip (11 kB)
  Preparing metadata (setup.py) ... done
Collecting bleach==3.1.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignme
nt3/assignment3//requirements.txt (line 3))
  Downloading bleach-3.1.0-py2.py3-none-any.whl.metadata (19 kB)
Collecting certifi==2019.6.16 (from -r /content/drive/My Drive/CS6353/Assignments/ass
ignment3/assignment3//requirements.txt (line 4))
  Downloading certifi-2019.6.16-py2.py3-none-any.whl.metadata (2.5 kB)
Collecting cycler==0.10.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 5))
  Downloading cycler-0.10.0-py2.py3-none-any.whl.metadata (722 bytes)
Collecting decorator==4.4.0 (from -r /content/drive/My Drive/CS6353/Assignments/assig
nment3/assignment3//requirements.txt (line 6))
  Downloading decorator-4.4.0-py2.py3-none-any.whl.metadata (3.7 kB)
Collecting defusedxml==0.6.0 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 7))
  Downloading defusedxml-0.6.0-py2.py3-none-any.whl.metadata (31 kB)
Collecting entrypoints==0.3 (from -r /content/drive/My Drive/CS6353/Assignments/assig
nment3/assignment3//requirements.txt (line 8))
  Downloading entrypoints-0.3-py2.py3-none-any.whl.metadata (1.4 kB)
Collecting future==0.17.1 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 9))
  Downloading future-0.17.1.tar.gz (829 kB)
     ──────────────────────────────────── 829.1/829.1 kB 33.8 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Collecting imageio==2.5.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 10))
  Downloading imageio-2.5.0-py3-none-any.whl.metadata (2.8 kB)
Collecting ipykernel==5.1.2 (from -r /content/drive/My Drive/CS6353/Assignments/assig
nment3/assignment3//requirements.txt (line 11))
  Downloading ipykernel-5.1.2-py3-none-any.whl.metadata (919 bytes)
```

```
Collecting ipython==7.8.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 12))
  Downloading ipython-7.8.0-py3-none-any.whl.metadata (4.3 kB)
Requirement already satisfied: ipython-genutils==0.2.0 in /usr/local/lib/python3.10/d
ist-packages (from -r /content/drive/My Drive/CS6353/Assignments/assignment3/assignme
nt3//requirements.txt (line 13)) (0.2.0)
Collecting ipywidgets==7.5.1 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 14))
  Downloading ipywidgets-7.5.1-py2.py3-none-any.whl.metadata (1.8 kB)
Collecting jedi==0.15.1 (from -r /content/drive/My Drive/CS6353/Assignments/assignmen
t3/assignment3//requirements.txt (line 15))
  Downloading jedi-0.15.1-py2.py3-none-any.whl.metadata (15 kB)
Collecting Jinja2==2.10.1 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 16))
  Downloading Jinja2-2.10.1-py2.py3-none-any.whl.metadata (2.2 kB)
Collecting jsonschema==3.0.2 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 17))
  Downloading jsonschema-3.0.2-py2.py3-none-any.whl.metadata (7.4 kB)
Collecting jupyter==1.0.0 (from -r /content/drive/My Drive/CS6353/Assignments/assignm
ent3/assignment3//requirements.txt (line 18))
  Downloading jupyter-1.0.0-py2.py3-none-any.whl.metadata (995 bytes)
Collecting jupyter-client==5.3.1 (from -r /content/drive/My Drive/CS6353/Assignments/
assignment3/assignment3//requirements.txt (line 19))
  Downloading jupyter_client-5.3.1-py2.py3-none-any.whl.metadata (3.6 kB)
Collecting jupyter-console==6.0.0 (from -r /content/drive/My Drive/CS6353/Assignment
s/assignment3/assignment3//requirements.txt (line 20))
  Downloading jupyter_console-6.0.0-py2.py3-none-any.whl.metadata (955 bytes)
Collecting jupyter-core==4.5.0 (from -r /content/drive/My Drive/CS6353/Assignments/as
signment3/assignment3//requirements.txt (line 21))
  Downloading jupyter_core-4.5.0-py2.py3-none-any.whl.metadata (884 bytes)
Collecting kiwisolver==1.1.0 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 22))
  Downloading kiwisolver-1.1.0.tar.gz (30 kB)
  Preparing metadata (setup.py) ... done
Collecting MarkupSafe==1.1.1 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 23))
  Downloading MarkupSafe-1.1.1.tar.gz (19 kB)
  Preparing metadata (setup.py) ... done
Collecting matplotlib==3.1.1 (from -r /content/drive/My Drive/CS6353/Assignments/assi
gnment3/assignment3//requirements.txt (line 24))
  Downloading matplotlib-3.1.1.tar.gz (37.8 MB)
     ──────────────────────────────────── 37.8/37.8 MB 27.7 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Requirement already satisfied: mistune==0.8.4 in /usr/local/lib/python3.10/dist-packa
ges (from -r /content/drive/My Drive/CS6353/Assignments/assignment3/assignment3//requ
irements.txt (line 25)) (0.8.4)
ERROR: Could not find a version that satisfies the requirement mkl-fft==1.0.6 (from v
ersions: 1.3.6, 1.3.8)
ERROR: No matching distribution found for mkl-fft==1.0.6
```

```python
In [ ]:  # As usual, a bit of setup
         from __future__ import print_function
         import time
         import numpy as np
         import matplotlib.pyplot as plt
         from cs6353.classifiers.fc_net import *
         from cs6353.data_utils import get_CIFAR10_data
         from cs6353.gradient_check import eval_numerical_gradient, eval_numerical_gradient_arr
         from cs6353.solver import Solver
```

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

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

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

```python
In [ ]:  # Load the (preprocessed) CIFAR10 data.
data = get_CIFAR10_data()
for k, v in list(data.items()):
  print(('%s: ' % k, v.shape))
```

```
('X_train: ', (49000, 3, 32, 32))
('y_train: ', (49000,))
('X_val: ', (1000, 3, 32, 32))
('y_val: ', (1000,))
('X_test: ', (1000, 3, 32, 32))
('y_test: ', (1000,))
```

# Affine layer: foward

Open the file `cs6353/layers.py` and implement the `affine_forward` function.

Once you are done you can test your implementaion by running the following:

```python
In [ ]:  # Test the affine_forward function

num_inputs = 2
input_shape = (4, 5, 6)
output_dim = 3

input_size = num_inputs * np.prod(input_shape)
weight_size = output_dim * np.prod(input_shape)

x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape), output_dim)
b = np.linspace(-0.3, 0.1, num=output_dim)

out, _ = affine_forward(x, w, b)
correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
                        [ 3.25553199,  3.5141327,   3.77273342]])

# Compare your output with ours. The error should be around e-9 or less.
print('Testing affine_forward function:')
print('difference: ', rel_error(out, correct_out))
```

```
Testing affine_forward function:
difference:  9.769849468192957e-10
```

# Affine layer: backward

Now implement the `affine_backward` function and test your implementation using numeric gradient checking.

```
In [ ]:  # Test the affine_backward function
         np.random.seed(231)
         x = np.random.randn(10, 2, 3)
         w = np.random.randn(6, 5)
         b = np.random.randn(5)
         dout = np.random.randn(10, 5)

         dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0], x, dout)
         dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0], w, dout)
         db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0], b, dout)

         _, cache = affine_forward(x, w, b)
         dx, dw, db = affine_backward(dout, cache)

         # The error should be around e-10 or less
         print('Testing affine_backward function:')
         print('dx error: ', rel_error(dx_num, dx))
         print('dw error: ', rel_error(dw_num, dw))
         print('db error: ', rel_error(db_num, db))
```

```
Testing affine_backward function:
dx error:  5.399100368651805e-11
dw error:  9.904211865398145e-11
db error:  2.4122867568119087e-11
```

# ReLU activation: forward

Implement the forward pass for the ReLU activation function in the `relu_forward` function and test your implementation using the following:

```
In [ ]:  # Test the relu_forward function

         x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)

         out, _ = relu_forward(x)
         correct_out = np.array([[ 0.,          0.,          0.,          0.,         ],
                                 [ 0.,          0.,          0.04545455,  0.13636364,],
                                 [ 0.22727273,  0.31818182,  0.40909091,  0.5,        ]])

         # Compare your output with ours. The error should be on the order of e-8
         print('Testing relu_forward function:')
         print('difference: ', rel_error(out, correct_out))
```

```
Testing relu_forward function:
difference:  4.999999798022158e-08
```

# ReLU activation: backward

Now implement the backward pass for the ReLU activation function in the `relu_backward` function and test your implementation using numeric gradient checking:

```
In [ ]:  np.random.seed(231)
         x = np.random.randn(10, 10)
         dout = np.random.randn(*x.shape)

         dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout)

         _, cache = relu_forward(x)
         dx = relu_backward(dout, cache)

         # The error should be on the order of e-12
         print('Testing relu_backward function:')
         print('dx error: ', rel_error(dx_num, dx))
```

```
Testing relu_backward function:
dx error:  3.2756349136310288e-12
```

# Inline Question 1:

We've only asked you to implement ReLU, but there are a number of different activation functions that one could use in neural networks, each with its pros and cons. In particular, an issue commonly seen with activation functions is getting zero (or close to zero) gradient flow during backpropagation. Which of the following activation functions have this problem? If you consider these functions in the one dimensional case, what types of input would lead to this behaviour?

1. Sigmoid
2. ReLU
3. Leaky ReLU

# Answer:

1. Sigmoid: The sigmoid function is prone to the vanishing gradient problem because its gradient diminishes near zero when the input values are very large or very small. For instance, input values like [-10000, 100000] would cause this issue as the function enters saturation, where the gradients become tiny.

2. ReLU: ReLU, compared to sigmoid, is less prone to vanishing gradients because it responds linearly to positive inputs. Its gradient is either 1 for positive inputs or 0 for negative inputs. However, when all inputs are negative, the output becomes zero, leading to no gradient flow, which is when we call it the "dying ReLU" problem. An example of this would be if the inputs were all negative, like [-12, -13, -14].

3. Leaky ReLU: To address the problem of dying RELU, Leaky ReLU introduces a small negative slope for negative input values, i.e., for x < 0, it outputs 0.0001x, and for x ≥ 0, it returns x. This modification helps mitigate the vanishing gradient problem. However, since the

function isn't continuous at x = 0, the gradient is undefined at that point. If this isn't handled properly in the implementation, an example input like [0, 0, 0] could result in zero gradients.

# "Sandwich" layers

There are some common patterns of layers that are frequently used in neural nets. For example, affine layers are frequently followed by a ReLU nonlinearity. To make these common patterns easy, we define several convenience layers in the file `cs6353/layer_utils.py`.

For now take a look at the `affine_relu_forward` and `affine_relu_backward` functions, and run the following to numerically gradient check the backward pass:

```
In [ ]:  from cs6353.layer_utils import affine_relu_forward, affine_relu_backward
         np.random.seed(231)
         x = np.random.randn(2, 3, 4)
         w = np.random.randn(12, 10)
         b = np.random.randn(10)
         dout = np.random.randn(2, 10)

         out, cache = affine_relu_forward(x, w, b)
         dx, dw, db = affine_relu_backward(dout, cache)

         dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w, b)[0], x, d
         dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w, d
         db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], b, d

         # Relative error should be around e-10 or less
         print('Testing affine_relu_forward and affine_relu_backward:')
         print('dx error: ', rel_error(dx_num, dx))
         print('dw error: ', rel_error(dw_num, dw))
         print('db error: ', rel_error(db_num, db))
```

```
Testing affine_relu_forward and affine_relu_backward:
dx error:  2.299579177309368e-11
dw error:  8.162011105764925e-11
db error:  7.826724021458994e-12
```

# Loss layers: Softmax and SVM

You implemented these loss functions in the last assignment, so we'll give them to you for free here. You should still make sure you understand how they work by looking at the implementations in `cs6353/layers.py`.

You can make sure that the implementations are correct by running the following:

```
In [ ]:  np.random.seed(231)
         num_classes, num_inputs = 10, 50
         x = 0.001 * np.random.randn(num_inputs, num_classes)
         y = np.random.randint(num_classes, size=num_inputs)
```

```
dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be around the or
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x, verbose=False)
loss, dx = softmax_loss(x, y)

# Test softmax_loss function. Loss should be close to 2.3 and dx error should be aroun
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing svm_loss:
loss:  8.999602749096233
dx error:  1.4021566006651672e-09

Testing softmax_loss:
loss:  2.302545844500738
dx error:  9.384673161989355e-09
```

# Two-layer network

In the previous assignment you implemented a two-layer neural network in a single monolithic class. Now that you have implemented modular versions of the necessary layers, you will reimplement the two layer network using these modular implementations.

Open the file `cs6353/classifiers/fc_net.py` and complete the implementation of the `TwoLayerNet` class. This class will serve as a model for the other networks you will implement in this assignment, so read through it to make sure you understand the API. You can run the cell below to test your implementation.

```
In [ ]:  np.random.seed(231)
         N, D, H, C = 3, 5, 50, 7
         X = np.random.randn(N, D)
         y = np.random.randint(C, size=N)

         std = 1e-3
         model = TwoLayerNet(input_dim=D, hidden_dim=H, num_classes=C, weight_scale=std)

         print('Testing initialization ... ')
         W1_std = abs(model.params['W1'].std() - std)
         b1 = model.params['b1']
         W2_std = abs(model.params['W2'].std() - std)
         b2 = model.params['b2']
         assert W1_std < std / 10, 'First layer weights do not seem right'
         assert np.all(b1 == 0), 'First layer biases do not seem right'
         assert W2_std < std / 10, 'Second layer weights do not seem right'
         assert np.all(b2 == 0), 'Second layer biases do not seem right'

         print('Testing test-time forward pass ... ')
         model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
```

```python
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
  [[11.53165108, 12.2917344,  13.05181771, 13.81190102, 14.57198434, 15.33206765,
    [12.05769098, 12.74614105, 13.43459113, 14.1230412,  14.81149128, 15.49994135,
    [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506,
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

print('Testing training loss (no regularization)')
y = np.asarray([0, 5, 1])
loss, grads = model.loss(X, y)
correct_loss = 3.4702243556
assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'

model.reg = 1.0
loss, grads = model.loss(X, y)
correct_loss = 26.5948426952
assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'

# Errors should be around e-7 or less
for reg in [0.0, 0.7]:
  print('Running numeric gradient check with reg = ', reg)
  model.reg = reg
  loss, grads = model.loss(X, y)

  for name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
    print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg =  0.0
W1 relative error: 1.83e-08
W2 relative error: 3.12e-10
b1 relative error: 9.83e-09
b2 relative error: 4.33e-10
Running numeric gradient check with reg =  0.7
W1 relative error: 2.53e-07
W2 relative error: 2.85e-08
b1 relative error: 1.56e-08
b2 relative error: 7.76e-10
```

# Solver

In the previous assignment, the logic for training models was coupled to the models themselves. Following a more modular design, for this assignment we have split the logic for training models into a separate class.

Open the file `cs6353/solver.py` and read through it to familiarize yourself with the API. After doing so, use a `Solver` instance to train a `TwoLayerNet` that achieves at least `50%` accuracy

on the validation set.

```
In [ ]:  model = TwoLayerNet()
         solver = None

         ################################################################################
         # TODO: Use a Solver instance to train a TwoLayerNet that achieves at least  #
         # 50% accuracy on the validation set.                                        #
         ################################################################################

         # lr_list = [0.0001, 0.0002, 0.0003, 0.0004]
         # reg_list = [0.00001, 0.000015, 0.00002, 0.000025]
         # layer_list = [25, 50, 100]
         # epoch_list = [5, 10, 15, 20]

         lr_list = [0.0001, 0.0002, 0.0003]
         reg_list = [0.00001, 0.000015, 0.00002, 0.000025]
         layer_list = [50, 100]
         epoch_list = [10, 15, 20]

         # lr_list = [0.0003]
         # reg_list = [0.000025]
         # layer_list = [100]
         # epoch_list = [20]

         best_val_acc = -1
         for lr in lr_list:
           for regularization in reg_list:
             for layer in layer_list:
               for epoch in epoch_list:
                 print('lr: {}, reg: {}, layer: {}, epoch: {}'.format(lr, regularization, layer
                 model = TwoLayerNet(hidden_dim=layer, reg=regularization)
                 candidate_solver = Solver(model, data, optim_config={'learning_rate':lr},
                                           lr_decay = 0.95, num_epochs=epoch, verbose=False)
                 candidate_solver.train()
                 if candidate_solver.best_val_acc > best_val_acc:
                   best_val_acc = candidate_solver.best_val_acc
                   solver = candidate_solver
                   print('lr: {}, reg: {}, layer: {}, epoch: {}, Val acc: {}'.format(lr, regula

         print('Best val acc: {}'.format(best_val_acc))
         ################################################################################
         #                            END OF YOUR CODE                                #
         ################################################################################
```

```
lr: 0.0001, reg: 1e-05, layer: 50, epoch: 10
lr: 0.0001, reg: 1e-05, layer: 50, epoch: 10, Val acc: 0.461
lr: 0.0001, reg: 1e-05, layer: 50, epoch: 15
lr: 0.0001, reg: 1e-05, layer: 50, epoch: 15, Val acc: 0.477
lr: 0.0001, reg: 1e-05, layer: 50, epoch: 20
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 10
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 15
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 15, Val acc: 0.487
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 20, Val acc: 0.492
lr: 0.0001, reg: 1.5e-05, layer: 50, epoch: 10
lr: 0.0001, reg: 1.5e-05, layer: 50, epoch: 15
lr: 0.0001, reg: 1.5e-05, layer: 50, epoch: 20
lr: 0.0001, reg: 1.5e-05, layer: 100, epoch: 10
lr: 0.0001, reg: 1.5e-05, layer: 100, epoch: 15
lr: 0.0001, reg: 1.5e-05, layer: 100, epoch: 20
lr: 0.0001, reg: 1.5e-05, layer: 100, epoch: 20, Val acc: 0.496
lr: 0.0001, reg: 2e-05, layer: 50, epoch: 10
lr: 0.0001, reg: 2e-05, layer: 50, epoch: 15
lr: 0.0001, reg: 2e-05, layer: 50, epoch: 20
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 10
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 15
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 20, Val acc: 0.499
lr: 0.0001, reg: 2.5e-05, layer: 50, epoch: 10
lr: 0.0001, reg: 2.5e-05, layer: 50, epoch: 15
lr: 0.0001, reg: 2.5e-05, layer: 50, epoch: 20
lr: 0.0001, reg: 2.5e-05, layer: 100, epoch: 10
lr: 0.0001, reg: 2.5e-05, layer: 100, epoch: 15
lr: 0.0001, reg: 2.5e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 1e-05, layer: 50, epoch: 10
lr: 0.0002, reg: 1e-05, layer: 50, epoch: 10, Val acc: 0.505
lr: 0.0002, reg: 1e-05, layer: 50, epoch: 15
lr: 0.0002, reg: 1e-05, layer: 50, epoch: 20
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 10
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 15
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 20, Val acc: 0.527
lr: 0.0002, reg: 1.5e-05, layer: 50, epoch: 10
lr: 0.0002, reg: 1.5e-05, layer: 50, epoch: 15
lr: 0.0002, reg: 1.5e-05, layer: 50, epoch: 20
lr: 0.0002, reg: 1.5e-05, layer: 100, epoch: 10
lr: 0.0002, reg: 1.5e-05, layer: 100, epoch: 15
lr: 0.0002, reg: 1.5e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 2e-05, layer: 50, epoch: 10
lr: 0.0002, reg: 2e-05, layer: 50, epoch: 15
lr: 0.0002, reg: 2e-05, layer: 50, epoch: 20
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 10
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 15
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 20, Val acc: 0.545
lr: 0.0002, reg: 2.5e-05, layer: 50, epoch: 10
lr: 0.0002, reg: 2.5e-05, layer: 50, epoch: 15
lr: 0.0002, reg: 2.5e-05, layer: 50, epoch: 20
lr: 0.0002, reg: 2.5e-05, layer: 100, epoch: 10
lr: 0.0002, reg: 2.5e-05, layer: 100, epoch: 15
lr: 0.0002, reg: 2.5e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 1e-05, layer: 50, epoch: 10
lr: 0.0003, reg: 1e-05, layer: 50, epoch: 15
lr: 0.0003, reg: 1e-05, layer: 50, epoch: 20
```

```
lr: 0.0003, reg: 1e-05, layer: 100, epoch: 10
lr: 0.0003, reg: 1e-05, layer: 100, epoch: 15
lr: 0.0003, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 1.5e-05, layer: 50, epoch: 10
lr: 0.0003, reg: 1.5e-05, layer: 50, epoch: 15
lr: 0.0003, reg: 1.5e-05, layer: 50, epoch: 20
lr: 0.0003, reg: 1.5e-05, layer: 100, epoch: 10
lr: 0.0003, reg: 1.5e-05, layer: 100, epoch: 15
lr: 0.0003, reg: 1.5e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 2e-05, layer: 50, epoch: 10
lr: 0.0003, reg: 2e-05, layer: 50, epoch: 15
lr: 0.0003, reg: 2e-05, layer: 50, epoch: 20
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 10
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 15
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 2.5e-05, layer: 50, epoch: 10
lr: 0.0003, reg: 2.5e-05, layer: 50, epoch: 15
lr: 0.0003, reg: 2.5e-05, layer: 50, epoch: 20
lr: 0.0003, reg: 2.5e-05, layer: 100, epoch: 10
lr: 0.0003, reg: 2.5e-05, layer: 100, epoch: 15
lr: 0.0003, reg: 2.5e-05, layer: 100, epoch: 20
Best val acc: 0.545
```
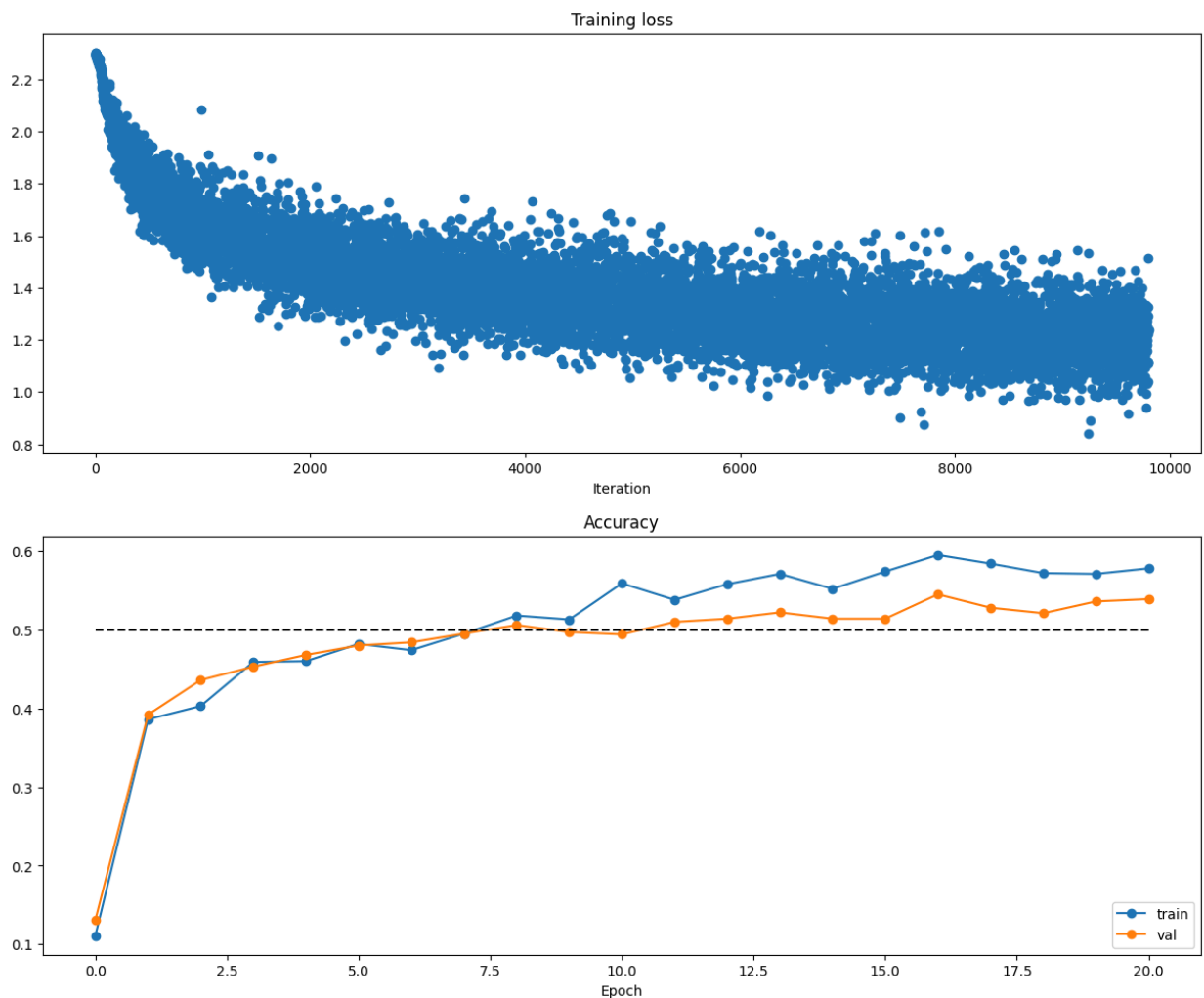
In [ ]:
```python
# Run this cell to visualize training loss and train / val accuracy

plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
plt.xlabel('Iteration')

plt.subplot(2, 1, 2)
plt.title('Accuracy')
plt.plot(solver.train_acc_history, '-o', label='train')
plt.plot(solver.val_acc_history, '-o', label='val')
plt.plot([0.5] * len(solver.val_acc_history), 'k--')
plt.xlabel('Epoch')
plt.legend(loc='lower right')
plt.gcf().set_size_inches(15, 12)
plt.show()
```

# Multilayer network

Next you will implement a fully-connected network with an arbitrary number of hidden layers.

Read through the `FullyConnectedNet` class in the file `cs6353/classifiers/fc_net.py` .

Implement the initialization, the forward pass, and the backward pass. For the moment don't worry about implementing batch/layer normalization; we will add those features soon.

## Initial loss and gradient check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. Do the initial losses seem reasonable?

For gradient checking, you should expect to see errors around 1e-7 or less.

```
In [ ]:  np.random.seed(231)
         N, D, H1, H2, C = 2, 15, 20, 30, 10
         X = np.random.randn(N, D)
         y = np.random.randint(C, size=(N,))
```

```python
for reg in [0, 3.14]:
  print('Running check with reg = ', reg)
  model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                            reg=reg, weight_scale=5e-2, dtype=np.float64)

  loss, grads = model.loss(X, y)
  print('Initial loss: ', loss)

  # Most of the errors should be on the order of e-7 or smaller.
  # NOTE: It is fine however to see an error for W2 on the order of e-5
  # for the check when reg = 0.0
  for name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
    print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Running check with reg =  0
Initial loss:  2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11
Running check with reg =  3.14
Initial loss:  7.052114776533016
W1 relative error: 3.90e-09
W2 relative error: 6.87e-08
W3 relative error: 2.13e-08
b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.57e-10
```

As another sanity check, make sure you can overfit a small dataset of 50 images. First we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the learning rate and initialization scale to overfit and achieve 100% training accuracy within 20 epochs.

```python
In [ ]:  # TODO: Use a three-layer Net to overfit 50 training examples by
         # tweaking just the learning rate and initialization scale.

         num_train = 50
         small_data = {
           'X_train': data['X_train'][:num_train],
           'y_train': data['y_train'][:num_train],
           'X_val': data['X_val'],
           'y_val': data['y_val'],
         }

         weight_scale = 1.5e-2
         learning_rate = 3e-3
         model = FullyConnectedNet([100, 100],
                     weight_scale=weight_scale, dtype=np.float64)
         solver = Solver(model, small_data,
                     print_every=10, num_epochs=20, batch_size=25,
                     update_rule='sgd',
                     optim_config={
                         'learning_rate': learning_rate,
                     }
```
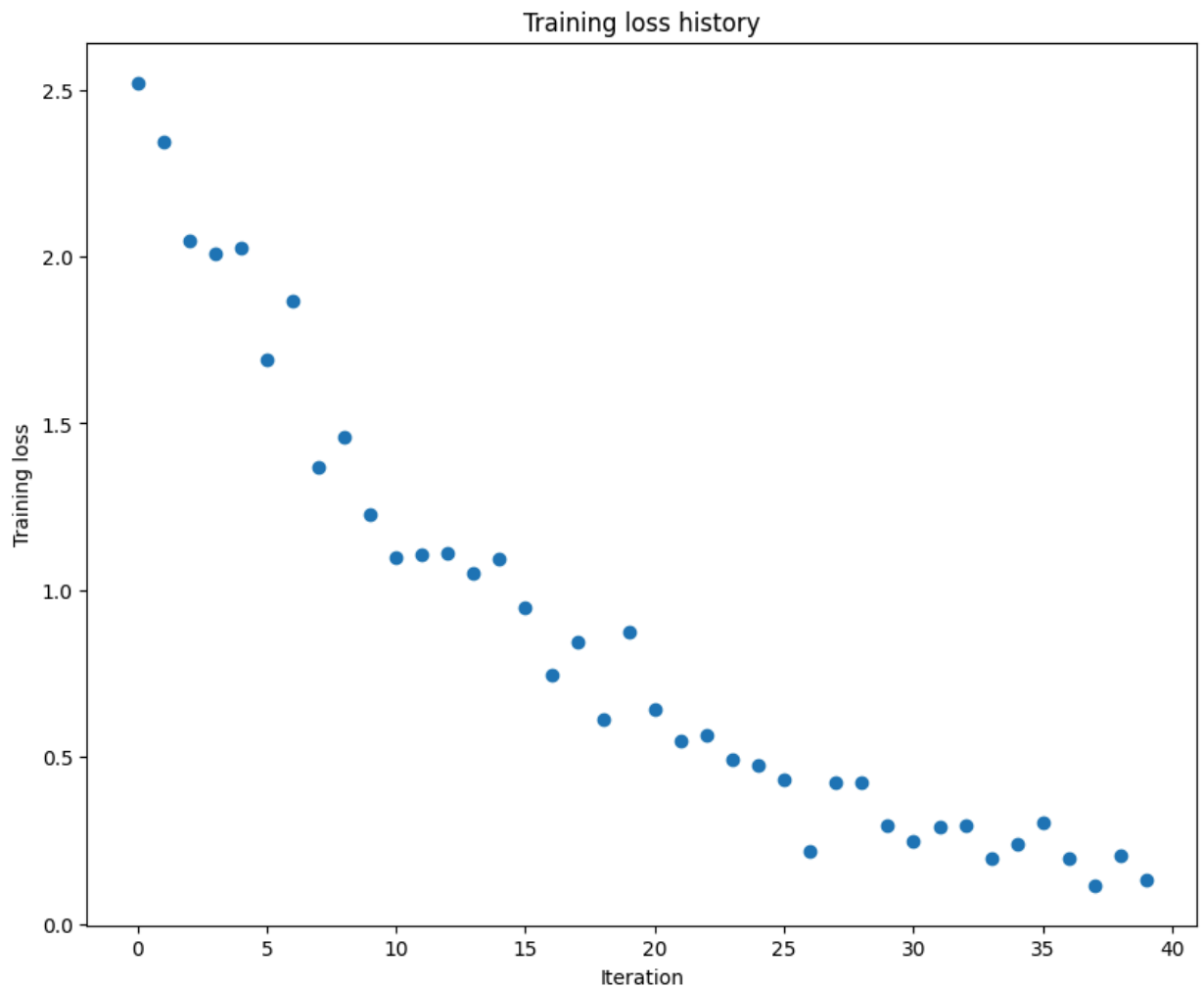
```
            )
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```

```
(Iteration 1 / 40) loss: 2.521152
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.121000
(Epoch 1 / 20) train acc: 0.240000; val_acc: 0.151000
(Epoch 2 / 20) train acc: 0.440000; val_acc: 0.153000
(Epoch 3 / 20) train acc: 0.460000; val_acc: 0.173000
(Epoch 4 / 20) train acc: 0.520000; val_acc: 0.168000
(Epoch 5 / 20) train acc: 0.700000; val_acc: 0.170000
(Iteration 11 / 40) loss: 1.097557
(Epoch 6 / 20) train acc: 0.840000; val_acc: 0.187000
(Epoch 7 / 20) train acc: 0.800000; val_acc: 0.173000
(Epoch 8 / 20) train acc: 0.880000; val_acc: 0.202000
(Epoch 9 / 20) train acc: 0.920000; val_acc: 0.192000
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.180000
(Iteration 21 / 40) loss: 0.641836
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.156000
(Epoch 12 / 20) train acc: 0.980000; val_acc: 0.176000
(Epoch 13 / 20) train acc: 0.980000; val_acc: 0.169000
(Epoch 14 / 20) train acc: 0.980000; val_acc: 0.188000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.196000
(Iteration 31 / 40) loss: 0.249410
(Epoch 16 / 20) train acc: 0.980000; val_acc: 0.193000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.189000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.192000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.193000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.187000
```

## Training loss history



Now try to use a five-layer network with 100 units on each layer to overfit 50 training examples. Again you will have to adjust the learning rate and weight initialization, but you should be able to achieve 100% training accuracy within 20 epochs.
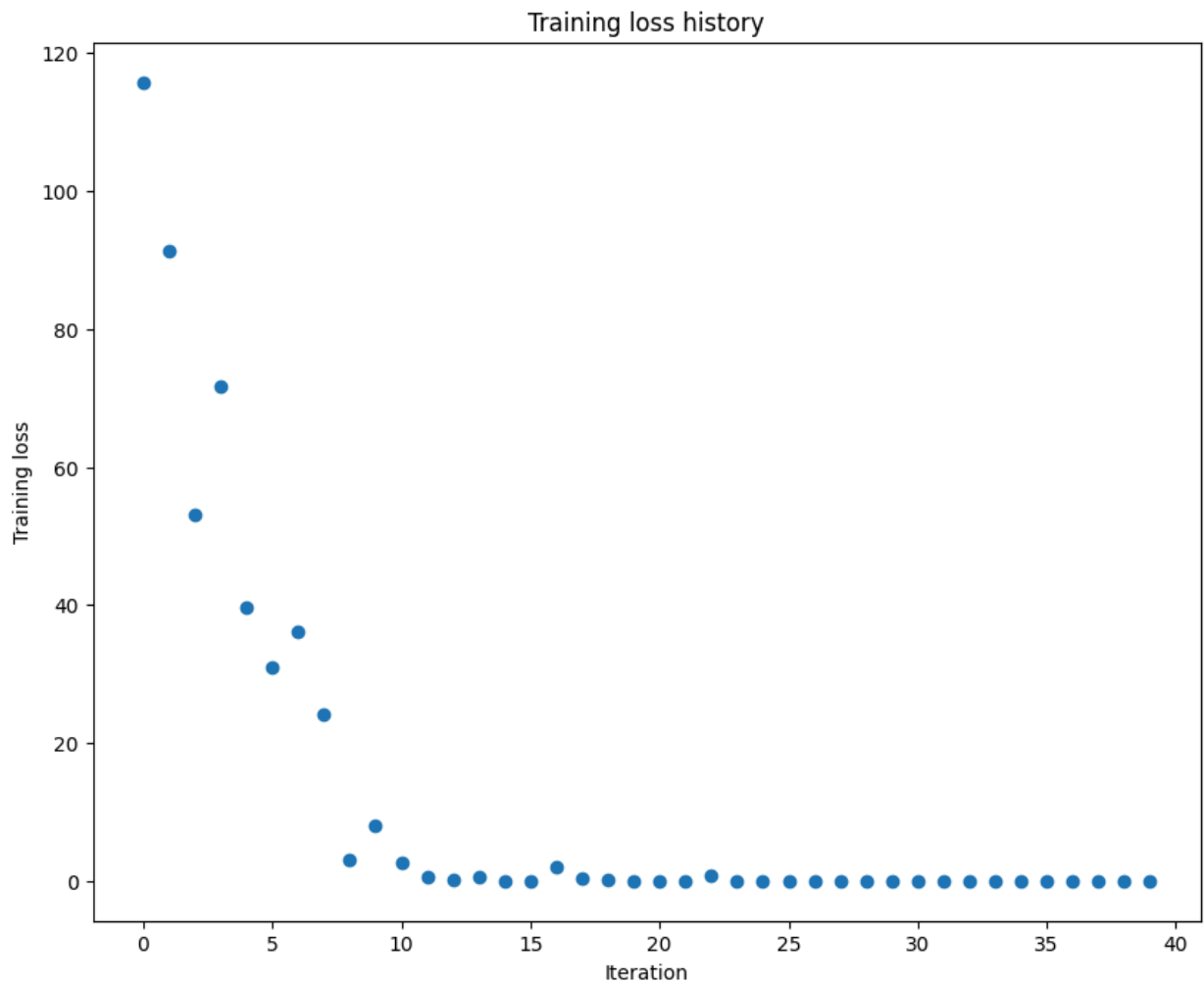
In [35]:
```python
# TODO: Use a five-layer Net to overfit 50 training examples by
# tweaking just the learning rate and initialization scale.

num_train = 50
small_data = {
  'X_train': data['X_train'][:num_train],
  'y_train': data['y_train'][:num_train],
  'X_val': data['X_val'],
  'y_val': data['y_val'],
}

# learning_rate = 2e-3
learning_rate = 1e-3
weight_scale = 1e-1
model = FullyConnectedNet([100, 100, 100, 100],
                weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                }
```

```python
        )
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```

```
(Iteration 1 / 40) loss: 115.774329
(Epoch 0 / 20) train acc: 0.140000; val_acc: 0.112000
(Epoch 1 / 20) train acc: 0.160000; val_acc: 0.126000
(Epoch 2 / 20) train acc: 0.320000; val_acc: 0.085000
(Epoch 3 / 20) train acc: 0.360000; val_acc: 0.143000
(Epoch 4 / 20) train acc: 0.680000; val_acc: 0.132000
(Epoch 5 / 20) train acc: 0.760000; val_acc: 0.133000
(Iteration 11 / 40) loss: 2.715416
(Epoch 6 / 20) train acc: 0.880000; val_acc: 0.132000
(Epoch 7 / 20) train acc: 0.940000; val_acc: 0.145000
(Epoch 8 / 20) train acc: 0.940000; val_acc: 0.144000
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.146000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.148000
(Iteration 21 / 40) loss: 0.000212
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.148000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.148000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.147000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.147000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.147000
(Iteration 31 / 40) loss: 0.000291
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.147000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.147000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.148000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.148000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.148000
```

Training loss history

## Inline Question 2:

Did you notice anything about the comparative difficulty of training the three-layer net vs training the five layer net? In particular, based on your experience, which network seemed more sensitive to the initialization scale? Why do you think that is the case?

## Answer:

Training a five-layer network is significantly more sensitive to weight initialization than a three-layer network due to its greater depth. In deep networks, small weight scales cause vanishing gradients, while large scales cause exploding gradients. As depth increases, the chance of encountering these issues rises, making it harder to find the right weight scale. The five-layer net was more sensetive to the weight initialization than the learning rate. Increasing the weight intialization causes exploding gradients quickly. While shallower networks also face sensitivity, their weight scale is easier to identify.

# Update rules

So far we have used vanilla stochastic gradient descent (SGD) as our update rule. More sophisticated update rules can make it easier to train deep networks. We will implement a few of the most commonly used update rules and compare them to vanilla SGD.

# SGD+Momentum

Stochastic gradient descent with momentum is a widely used update rule that tends to make deep networks converge faster than vanilla stochastic gradient descent.

Open the file `cs6353/optim.py` and read the documentation at the top of the file to make sure you understand the API. Implement the SGD+momentum update rule in the function `sgd_momentum` and run the following to check your implementation. You should see errors less than e-8.

```
In [36]:   from cs6353.optim import sgd_momentum


           N, D = 4, 5
           w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
           dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
           v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

           config = {'learning_rate': 1e-3, 'velocity': v}
           next_w, _ = sgd_momentum(w, dw, config=config)

           expected_next_w = np.asarray([
             [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
             [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
             [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
             [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096     ]])
           expected_velocity = np.asarray([
             [ 0.5406,      0.55475789,  0.56891579, 0.58307368,  0.59723158],
             [ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
             [ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
             [ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096     ]])

           # Should see relative errors around e-8 or less
           print('next_w error: ', rel_error(next_w, expected_next_w))
           print('velocity error: ', rel_error(expected_velocity, config['velocity']))
```

```
next_w error:   8.882347033505819e-09
velocity error:   4.269287743278663e-09
```

Once you have done so, run the following to train a six-layer network with both SGD and SGD+momentum. You should see the SGD+momentum update rule converge faster.

```
In [37]:   num_train = 4000
           small_data = {
             'X_train': data['X_train'][:num_train],
             'y_train': data['y_train'][:num_train],
             'X_val': data['X_val'],
             'y_val': data['y_val'],
           }
```

```python
solvers = {}

for update_rule in ['sgd', 'sgd_momentum']:
  print('running with ', update_rule)
  model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

  solver = Solver(model, small_data,
                  num_epochs=5, batch_size=100,
                  update_rule=update_rule,
                  optim_config={
                      'learning_rate': 1e-2,
                  },
                  verbose=True)
  solvers[update_rule] = solver
  solver.train()
  print()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)
plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in list(solvers.items()):
  plt.subplot(3, 1, 1)
  plt.plot(solver.loss_history, 'o', label=update_rule)

  plt.subplot(3, 1, 2)
  plt.plot(solver.train_acc_history, '-o', label=update_rule)

  plt.subplot(3, 1, 3)
  plt.plot(solver.val_acc_history, '-o', label=update_rule)

for i in [1, 2, 3]:
  plt.subplot(3, 1, i)
  plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```
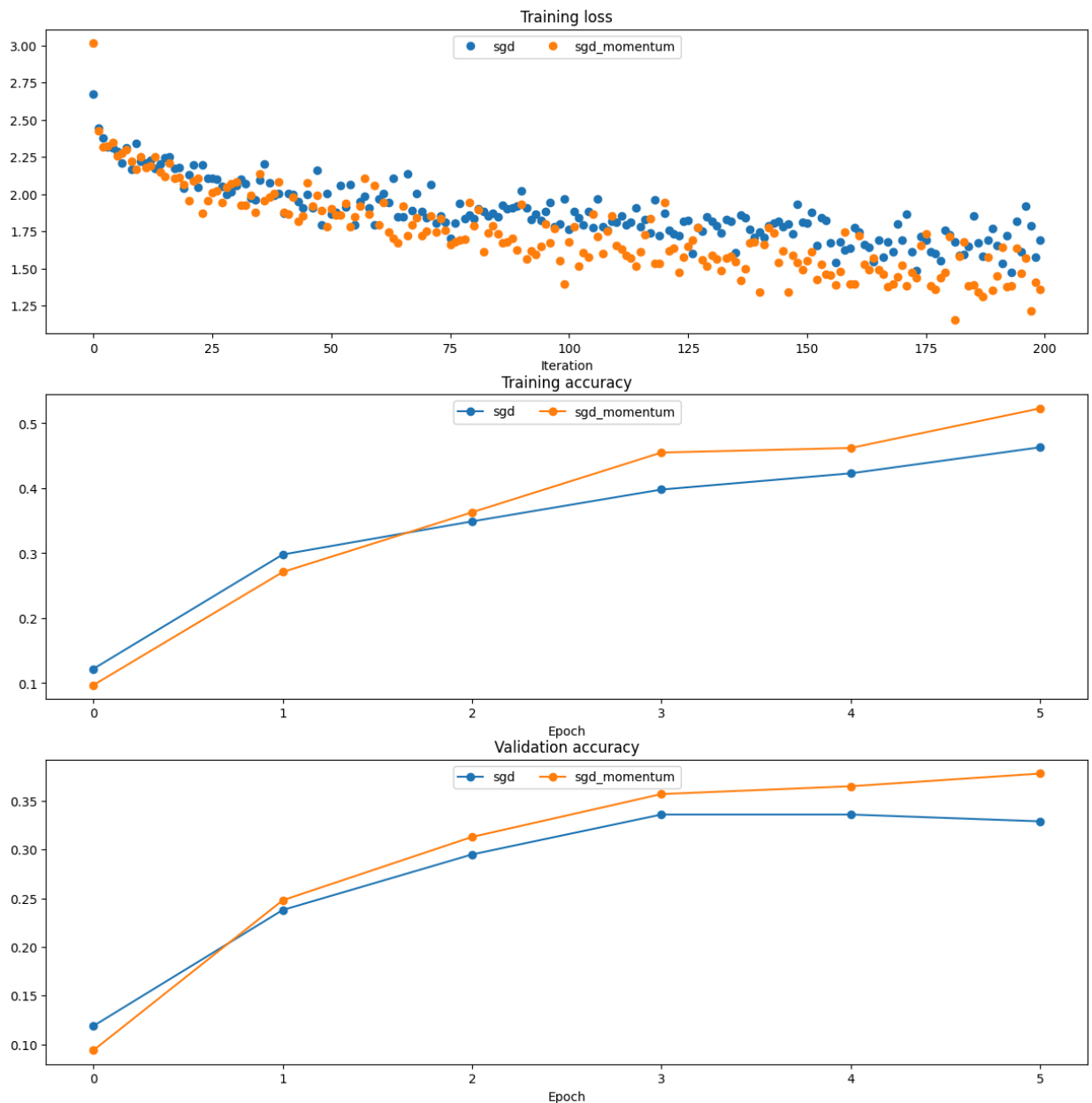
```
running with  sgd
(Iteration 1 / 200) loss: 2.675409
(Epoch 0 / 5) train acc: 0.122000; val_acc: 0.119000
(Iteration 11 / 200) loss: 2.221112
(Iteration 21 / 200) loss: 2.130692
(Iteration 31 / 200) loss: 2.058671
(Epoch 1 / 5) train acc: 0.298000; val_acc: 0.238000
(Iteration 41 / 200) loss: 1.869888
(Iteration 51 / 200) loss: 1.866691
(Iteration 61 / 200) loss: 1.967924
(Iteration 71 / 200) loss: 1.840155
(Epoch 2 / 5) train acc: 0.349000; val_acc: 0.295000
(Iteration 81 / 200) loss: 1.832665
(Iteration 91 / 200) loss: 2.024011
(Iteration 101 / 200) loss: 1.764401
(Iteration 111 / 200) loss: 1.813284
(Epoch 3 / 5) train acc: 0.398000; val_acc: 0.336000
(Iteration 121 / 200) loss: 1.869099
(Iteration 131 / 200) loss: 1.818794
(Iteration 141 / 200) loss: 1.745466
(Iteration 151 / 200) loss: 1.803553
(Epoch 4 / 5) train acc: 0.423000; val_acc: 0.336000
(Iteration 161 / 200) loss: 1.773484
(Iteration 171 / 200) loss: 1.688030
(Iteration 181 / 200) loss: 1.725408
(Iteration 191 / 200) loss: 1.651351
(Epoch 5 / 5) train acc: 0.463000; val_acc: 0.329000

running with  sgd_momentum
(Iteration 1 / 200) loss: 3.015264
(Epoch 0 / 5) train acc: 0.097000; val_acc: 0.094000
(Iteration 11 / 200) loss: 2.253214
(Iteration 21 / 200) loss: 1.954631
(Iteration 31 / 200) loss: 2.079628
(Epoch 1 / 5) train acc: 0.271000; val_acc: 0.248000
(Iteration 41 / 200) loss: 1.876052
(Iteration 51 / 200) loss: 1.903648
(Iteration 61 / 200) loss: 1.790516
(Iteration 71 / 200) loss: 1.751532
(Epoch 2 / 5) train acc: 0.363000; val_acc: 0.313000
(Iteration 81 / 200) loss: 1.787554
(Iteration 91 / 200) loss: 1.932285
(Iteration 101 / 200) loss: 1.680905
(Iteration 111 / 200) loss: 1.655849
(Epoch 3 / 5) train acc: 0.455000; val_acc: 0.357000
(Iteration 121 / 200) loss: 1.944164
(Iteration 131 / 200) loss: 1.587576
(Iteration 141 / 200) loss: 1.340575
(Iteration 151 / 200) loss: 1.554041
(Epoch 4 / 5) train acc: 0.462000; val_acc: 0.365000
(Iteration 161 / 200) loss: 1.393865
(Iteration 171 / 200) loss: 1.520558
(Iteration 181 / 200) loss: 1.716010
(Iteration 191 / 200) loss: 1.446959
(Epoch 5 / 5) train acc: 0.523000; val_acc: 0.378000
```

# Train a good model!

Train the best fully-connected model that you can on CIFAR-10, storing your best model in the `best_model` variable. We require you to get at least 50% accuracy on the validation set using a fully-connected net.

If you are careful it should be possible to get accuracies above 55%, but we don't require it for this part and won't assign extra credit for doing so. Later in the assignment we will ask you to train the best convolutional network that you can on CIFAR-10, and we would prefer that you spend your effort working on convolutional nets rather than fully-connected nets.

You might find it useful to complete the `BatchNormalization.ipynb` notebook before completing this part, since those techniques can help you train powerful models.

In [42]:
```python
best_model = None
################################################################################
# TODO: Train the best FullyConnectedNet that you can on CIFAR-10. You might   #
# find batch/layer normalization useful. Store your best model in  #
# the best_model variable.                                                     #
################################################################################
# lr_list = [0.0001, 0.0002, 0.0003, 0.0004]
# reg_list = [0.00001, 0.000015, 0.00002, 0.000025]
# layer_list = [25, 50, 100]
# epoch_list = [5, 10, 15, 20]

lr_list = [0.0001, 0.0002, 0.0003]
reg_list = [0.00001, 0.00002]
layer_list = [100]
epoch_list = [20, 25]

# lr_list = [0.0003]
# reg_list = [0.00002]
# layer_list = [100]
# epoch_list = [20]

best_val_acc = -1
for lr in lr_list:
  for regularization in reg_list:
    for layer in layer_list:
      for epoch in epoch_list:
        print('lr: {}, reg: {}, layer: {}, epoch: {}'.format(lr, regularization, layer
        model = TwoLayerNet(hidden_dim=layer, reg=regularization)
        candidate_solver = Solver(model, data, update_rule='sgd', optim_config={'learn
                                  lr_decay = 0.95, num_epochs=epoch, batch_size=100, \
        candidate_solver.train()
        if candidate_solver.best_val_acc > best_val_acc:
          best_val_acc = candidate_solver.best_val_acc
          best_model = model
          print('lr: {}, reg: {}, layer: {}, epoch: {}, Val acc: {}'.format(lr, regula

print('Best val acc: {}'.format(best_val_acc))
################################################################################
#                              END OF YOUR CODE                                #
################################################################################
```

```
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 20, Val acc: 0.482
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 25
lr: 0.0001, reg: 1e-05, layer: 100, epoch: 25, Val acc: 0.5
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 25
lr: 0.0001, reg: 2e-05, layer: 100, epoch: 25, Val acc: 0.504
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 20, Val acc: 0.525
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 25
lr: 0.0002, reg: 1e-05, layer: 100, epoch: 25, Val acc: 0.533
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0002, reg: 2e-05, layer: 100, epoch: 25
lr: 0.0003, reg: 1e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 1e-05, layer: 100, epoch: 25
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 20
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 25
lr: 0.0003, reg: 2e-05, layer: 100, epoch: 25, Val acc: 0.537
Best val acc: 0.537
```

# Test your model!

Run your best model on the validation and test sets. You should achieve above 50% accuracy on the validation set.

In [43]:
```python
y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)
print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())
print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())
```

```
Validation set accuracy:  0.537
Test set accuracy:  0.527
```