

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 for each model 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:

```
def rel_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:

```
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 Dropout as a regularizer and Batch/Layer Normalization as a tool to more efficiently optimize deep networks.

Acknowledgement: This exercise is adapted from [Stanford CS231n](#).

```
In [15]: # As usual, a bit of setup
from future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from libs.classifiers.fc_net import *
from libs.data_utils import get_cifar10_data
from libs.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from libs.solver import Solver

%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.max(abs(x) + np.abs(y)))

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

```
In [16]: # 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': (10000, 3, 32, 32))
('y_test': (1000,))
```

Affine layer: foward

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

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

```
In [17]: # Test the affine_forward function
np.random.seed(231)
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(-0.3, 0.3, num=weight_size).reshape(np.prod(input_shape), output_dim)
b = np.linspace(-0.2, 0.1, num=output_dim)

_, cache = affine_forward(x, w, b)
out, dx, dw, db = np.array([[ 1.49834967, 1.70660132, 1.91485297],
                             [ 3.2553199, 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 backward function:
difference: 9.76894468192957e-10
```

Affine layer: backward

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

```
In [18]: # 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: 5.399100366651805e-11
dw error: 9.304211865398145e-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 [19]: # 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],
                         [ 0.27272727, 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: 4.99999798022159e-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 [20]: 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
```

"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 `libs/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 [21]: from libs.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, dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], b, dout)

# 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: 8.16201105764925e-11
db error: 7.826724021458994e-12
```

Loss layers: Softmax

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 `libs/layers.py`.

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

```
In [22]: 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: 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 around e-8
print('Testing softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

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 `libs/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 [23]: np.random.seed(231)
N, D, H1, H2, 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)
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, 16.09215096],
     [12.59769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49994135, 16.18839143],
     [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506, 16.2846319 ]])
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, h=1e-5)
        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: 1.53e-09
b2 relative error: 4.13e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.53e-07
W2 relative error: 1.85e-08
b1 relative error: 1.56e-09
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 `libs/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 **58%** accuracy on the validation set.

```
In [24]: # X_val: (1000, 3, 32, 32)
# X_train: (49000, 3, 32, 32)
# X_test: (1000, 3, 32, 32)
# y_val: (1000,)
# y_train: (49000,)
# y_test: (1000,)

# model = TwoLayerNet()
# solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least 58% accuracy
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

model = TwoLayerNet()
solver = Solver(model, data, num_epochs=10, batch_size=100, update_rule='sgd', optim_config={'learning_rate': 1e-3}, solver=train())

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

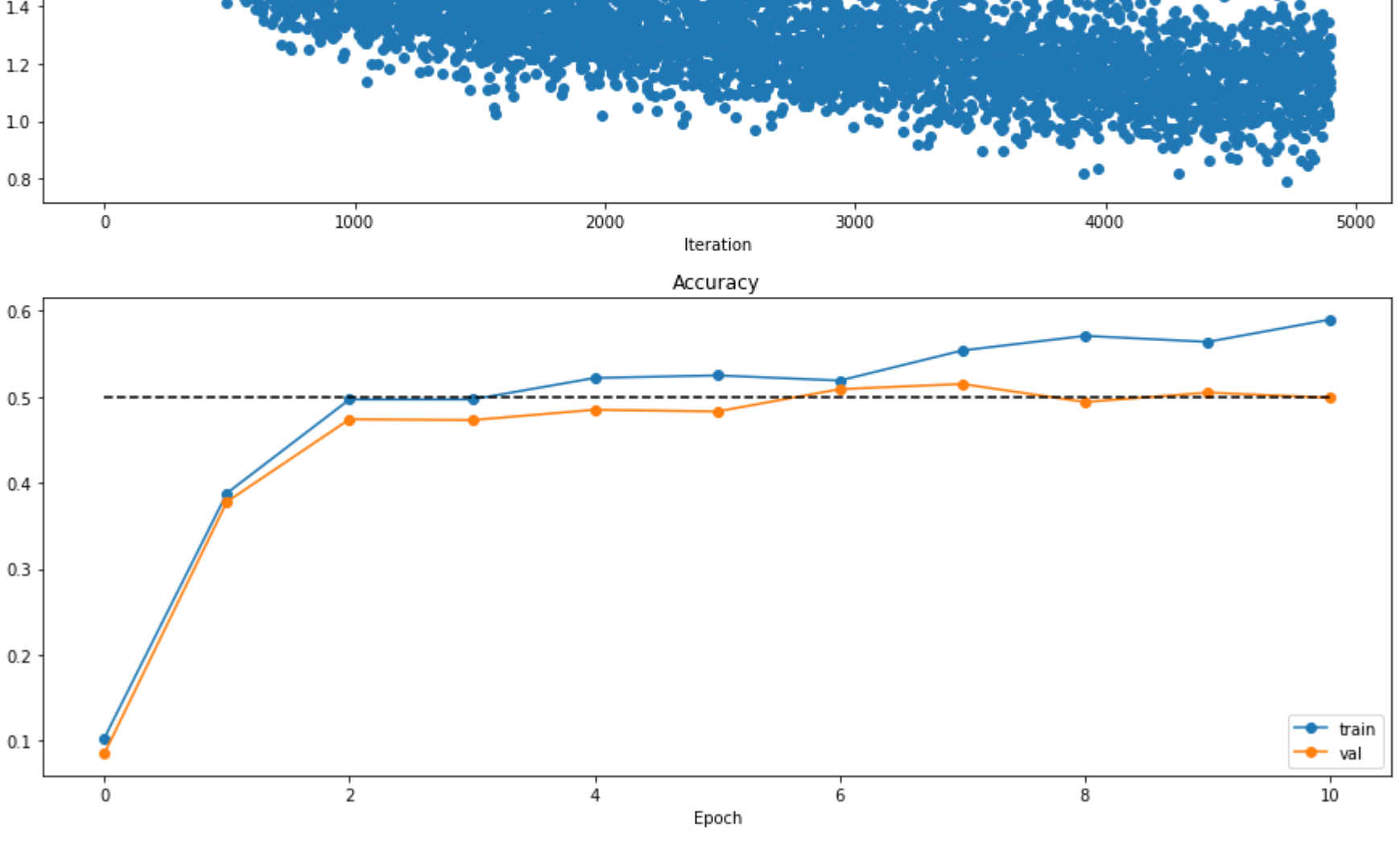
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####

(iteration 1 / 4900) loss: 2.305171
(iteration 0 / 10) train acc: 0.103000; val_acc: 0.085000
(iteration 201 / 4900) loss: 1.932779
(iteration 401 / 4900) loss: 1.538737
(epoch 1 / 10) train acc: 0.388000; val_acc: 0.378000
(iteration 601 / 4900) loss: 1.701701
(iteration 801 / 4900) loss: 1.699713
(epoch 2 / 10) train acc: 0.497000; val_acc: 0.474000
(iteration 1001 / 4900) loss: 1.424626
(iteration 1201 / 4900) loss: 1.666793
(iteration 1401 / 4900) loss: 1.211955
(epoch 3 / 10) train acc: 0.497000; val_acc: 0.473000
(iteration 1601 / 4900) loss: 1.281926
(iteration 1801 / 4900) loss: 1.357318
(epoch 4 / 10) train acc: 0.520000; val_acc: 0.485000
(iteration 2001 / 4900) loss: 1.311724
(iteration 2201 / 4900) loss: 1.228021
(iteration 2401 / 4900) loss: 1.36943
(epoch 5 / 10) train acc: 0.525000; val_acc: 0.483000
(iteration 2601 / 4900) loss: 1.404789
(iteration 2801 / 4900) loss: 1.238467
(epoch 6 / 10) train acc: 0.519000; val_acc: 0.509000
(iteration 3001 / 4900) loss: 1.216040
(iteration 3201 / 4900) loss: 1.295010
(iteration 3401 / 4900) loss: 1.426607
(epoch 7 / 10) train acc: 0.554000; val_acc: 0.515000
(iteration 3601 / 4900) loss: 1.121254
(iteration 3801 / 4900) loss: 1.124680
(epoch 8 / 10) train acc: 0.571000; val_acc: 0.494000
(iteration 4001 / 4900) loss: 1.221882
(iteration 4201 / 4900) loss: 1.237958
(iteration 4401 / 4900) loss: 1.325068
(epoch 9 / 10) train acc: 0.564000; val_acc: 0.505000
(iteration 4601 / 4900) loss: 1.339394
(iteration 4801 / 4900) loss: 1.071422
(epoch 10 / 10) train acc: 0.590000; val_acc: 0.499000
```

```
In [25]: # 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), 'x--')
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 `libs/classifiers/fc_net.py`.

Implement the initialization, the forward pass, and the backward pass. For the moment don't worry about implementing dropout or 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 [26]: 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,))

for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              update_rule='sgd', 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: 7.0521147633016
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 **weight initialization scale** to overfit and achieve 100% training accuracy within 20 epochs.

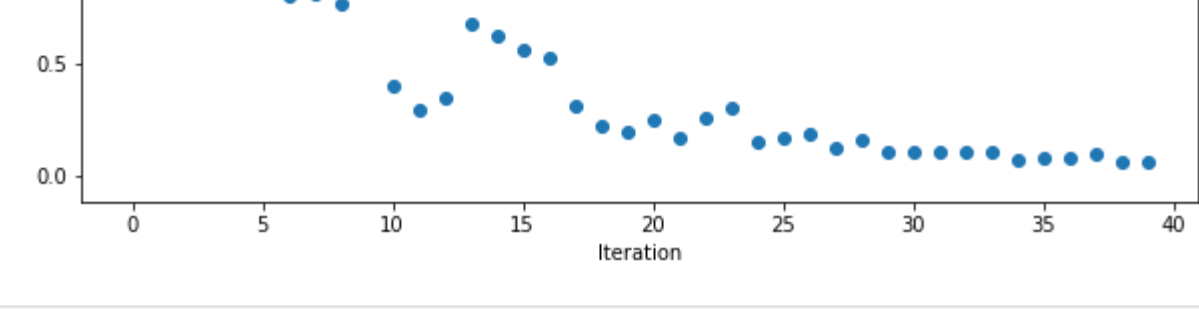
```
In [27]: # TODO: Use a three-layer Net to overfit 50 training examples by
# tweaking just the learning rate and initialization scale.

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 = 2e-2 # Experiment with this!
learning_rate = 3e-3 # Experiment with this!
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: 3.604568
(epoch 0 / 20) train acc: 0.180000; val_acc: 0.190000
(epoch 1 / 20) train acc: 0.360000; val_acc: 0.137000
(epoch 2 / 20) train acc: 0.600000; val_acc: 0.155000
(epoch 3 / 20) train acc: 0.720000; val_acc: 0.170000
(epoch 4 / 20) train acc: 0.800000; val_acc: 0.171000
(epoch 5 / 20) train acc: 0.880000; val_acc: 0.182000
(iteration 11 / 40) loss: 0.402117
(epoch 6 / 20) train acc: 0.880000; val_acc: 0.194000
(epoch 7 / 20) train acc: 0.900000; val_acc: 0.183000
(epoch 8 / 20) train acc: 0.940000; val_acc: 0.203000
(epoch 9 / 20) train acc: 0.980000; val_acc: 0.195000
(epoch 10 / 20) train acc: 0.980000; val_acc: 0.185000
(iteration 21 / 40) loss: 0.245577
(epoch 11 / 20) train acc: 1.000000; val_acc: 0.201000
(epoch 12 / 20) train acc: 1.000000; val_acc: 0.190000
(epoch 13 / 20) train acc: 1.000000; val_acc: 0.191000
(epoch 14 / 20) train acc: 1.000000; val_acc: 0.190000
(epoch 15 / 20) train acc: 1.000000; val_acc: 0.187000
(iteration 31 / 40) loss: 0.107372
(epoch 16 / 20) train acc: 1.000000; val_acc: 0.195000
(epoch 17 / 20) train acc: 1.000000; val_acc: 0.196000
(epoch 18 / 20) train acc: 1.000000; val_acc: 0.193000
(epoch 19 / 20) train acc: 1.000000; val_acc: 0.191000
(epoch 20 / 20) train acc: 1.000000; val_acc: 0.190000
```



```
In [ ]:
```