

Problem Set 3: Neural Networks

This assignment requires a working Python Notebook installation, which you should already have. If not, please refer to the instructions in the course resources.

The programming part is adapted from [Stanford CS231n](#).

In part 2 (programming) of this assignment, you DO NOT need to make any modification code in this IPython Notebook. Instead you will implement your own simple neural network in the mlp.py file. Please attach your written solutions for part 1 and part 3 in this IPython Notebook.

Total: 80 points.

[30pts] Problem 1: Backprop in a simple MLP

This problem asks you to derive all the steps of the backpropagation algorithm for a simple classification network. Consider a fully-connected neural network, also known as a multi-layer perceptron (MLP), with a single hidden layer and a one-node output layer. The hidden and output nodes use an elementwise sigmoid activation function and the loss layer uses cross-entropy loss:

$$f(z) = \frac{1}{1 + \exp(-z)}$$
$$L(y, \hat{y}) = -y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})$$

The computation graph for an example network is shown below. Note that it has an equal number of nodes in the input and hidden layer (3 each), but, in general, they need not be equal. Also, to make the application of backprop easier, we show the computation graph which shows the dot product and activation functions as their own nodes, rather than the usual graph showing a single node for both.



The forward and backward computation are given below. NOTE: We assume no regularization, so you can omit the terms involving Ω .

The forward step is:

And the backward step is:

Write down each step of the backward pass explicitly for all layers, i.e. for the loss and $k = 2, 1$, compute all gradients above, expressing them as a function of variables x, y, h, W, b . We start by giving an example. Note that \odot stands for element-wise multiplication.

$$\nabla_{y^{(2)}} L(y, \hat{y}) = \nabla_{y^{(2)}} [-y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})] = \frac{\partial}{\partial \hat{y}} [-y \ln(\hat{y}) - (1 - y) \ln(1 - \hat{y})] \cdot \frac{\partial \hat{y}}{\partial y^{(2)}} = \frac{y - \hat{y}}{(1 - \hat{y})^2} \cdot \frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2}$$

Next, please derive the following.

Hint: you should substitute the updated values for the gradient g in each step and simplify as much as possible.

Useful information about vectorized chain rule and backpropagation:
If you are struggling with computing the vectorized version of chain rule for the backpropagation question in problem set 4, you may find this example helpful: <https://web.stanford.edu/ics202d/lectures/gradient-roles.pdf>. It also contains some helpful shortcuts for computing gradients.

[5pts] Q1.1: $\nabla_{y^{(2)}} J$

$$\nabla_{y^{(2)}} J = \frac{\partial J}{\partial y^{(2)}} = \frac{\partial J}{\partial W^{(2)}} \frac{\partial W^{(2)}}{\partial y^{(2)}} = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} f'(y^{(2)}) = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right]$$

[5pts] Q1.2: $\nabla_{y^{(1)}} J$

$$\nabla_{y^{(1)}} J = \frac{\partial J}{\partial y^{(1)}} \frac{\partial y^{(1)}}{\partial W^{(1)}} = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] \nabla_{W^{(2)}} [H^{(1)} W^{(2)} + b^{(2)}] = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] H^{(1)}$$

[5pts] Q1.4: $\nabla_{y^{(1)}} J$

$$\nabla_{y^{(1)}} J = \frac{\partial J}{\partial y^{(1)}} \frac{\partial y^{(1)}}{\partial W^{(1)}} = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] \nabla_{W^{(1)}} [H^{(1)} W^{(2)} + b^{(2)}] = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] W^{(2)}$$

[5pts] Q1.5: $\nabla_{y^{(1)}} J, \nabla_{W^{(1)}} J, \nabla_{b^{(1)}} J$

$$\nabla_{y^{(1)}} J = \frac{\partial J}{\partial y^{(1)}} \frac{\partial y^{(1)}}{\partial W^{(1)}} = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] W^{(2)}$$
$$\nabla_{W^{(1)}} J = \frac{\partial J}{\partial W^{(1)}} \frac{\partial y^{(1)}}{\partial W^{(1)}} = \frac{W^{(2)} - y}{(1 - W^{(2)})^2} \left[\frac{e^{-y^{(2)}}}{(e^{-y^{(2)}} + 1)^2} \right] W^{(2)}$$

[5pts] Q1.6 Briefly, explain how the computational speed of backpropagation would be affected if it did not include a forward pass

The forward pass is necessary to perform backpropagation. Without the forward pass, backpropagation wouldn't work.

[30pts] Problem 2 (Programming): Implementing a simple MLP

In this problem we will develop a neural network with fully-connected layers, or Multi-Layer Perceptron (MLP). We will use it in classification tasks.

In the current directory, you can find a file `mlp.py`, which contains the definition for class `TwoLayerMLP`. As the name suggests, it implements a 2-layer MLP or MLP with 1 hidden layer. You will implement your code in the same file, and call the member functions in this notebook. Below is some initialization. The `autoreload` command makes sure that `mlp.py` is periodically reloaded.

```
In [1]: # setup
import numpy as np
import matplotlib.pyplot as plt
from mlp import TwoLayerMLP

#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.maximum(1e-8, np.abs(x) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

Next we initialize a toy model and some toy data, the task is to classify five 4-d vectors.

```
In [2]: # Create a small net and some toy data to check your implementations.
# Note that we set the random seed for repeatable experiments.
input_size = 4
hidden_size = 10
num_classes = 3
num_inputs = 5

def init_toy_model(activ, std=1e-1):
    np.random.seed(0)
    return TwoLayerMLP(input_size, hidden_size, num_classes, std=std, activation=activ)

def init_toy_data():
    np.random.seed(1)
    y = 10 * np.random.randn(num_inputs, input_size)
    x = np.array([0, 1, 2, 2, 1])
    return x, y

x, y = init_toy_data()
print('X = ', X)
print('y = ', y)

X = [[ 16.24345364 -6.11756414 -5.28171752 -10.72686822]
 [ 0.64016219 -23.01336597 17.44811764 -7.61206901]
 [ 3.19030996 -2.49370375 14.62107937 -20.60140709]
 [-3.22417204 -3.84054355 11.33769442 -10.99891267]
 [-1.72426208 -8.77858418 0.42213747 5.82815214]]

y = [0 1 2 2 1]
```

[5pts] Q2.1 Forward pass: Sigmoid

Our 2-layer MLP uses a softmax output layer (note: this means that you don't need to apply a sigmoid on the output) and the multiclass cross-entropy loss to perform classification.

Softmax function:

For class j :

$$P(y_{(j)}|x) = \frac{\exp(z_j)}{\sum_{c=1}^C \exp(z_c)}$$

Where C is the number of classes and z is class-wise output of the network.

Multiclass cross-entropy loss function:

$$J = \frac{1}{m} \sum_{i=1}^m \sum_{c=1}^C [-y_{(c)} \log(P(y_{(c)}|x^{(i)}))]$$

$y_{(c)} = 1$ for the ground truth class and 0 otherwise.

m is the number of inputs in a batch and C is the number of classes.

Please take a look at method `TwoLayerMLP.loss` in the file `mlp.py`. This function takes in the data and weight parameters, and computes the class scores (aka logits), the loss J , and the gradients on the parameters.

- Complete the implementation of forward pass (up to the computation of `scores`) for the sigmoid activation: $\sigma(x) = \frac{1}{1 + \exp(-x)}$.

Note 1: Softmax cross entropy loss involves the `log-sum-exp` operation. This can result in numerical underflow/overflow. Read about the solution in the link, and try to understand the calculation of `loss` in the code.

Note 2: You're strongly encouraged to implement in a vectorized way and avoid using slower `for` loops. Note that most numpy functions support vector inputs.

Check the correctness of your forward pass below. The difference should be very small (<1e-6).

```
In [3]: net = init_toy_model('sigmoid', std=1e-1)
loss, _ = net.loss(X, y, reg=0.1)
correct_loss = 1.182248
print(loss)
print('Difference between your loss and correct loss:')
print(np.sum(np.abs(loss - correct_loss)))

1.1822479803941373
Difference between your loss and correct loss:
1.9605862711102873e-08

In [ ]:
```

[10pts] Q2.2 Backward pass: Sigmoid

- For sigmoid activation, complete the computation of `grads`, which stores the gradient of the loss with respect to the variables `W1`, `b1`, `W2`, and `b2`.

Now debug your backward pass using a numeric gradient check. Again, the differences should be very small.

```
In [4]: # Use numeric gradient checking to check your implementation of the backward pass.
# If your implementation is correct, the difference between the numeric and analytic gradients should be less than 1e-8 for each of W1, W2, b1, and b2.
from util import eval_numerical_gradient

loss, grads = net.loss(X, y, reg=0.1)

# these should all be very small
for param_name in grads:
    f = lambda W: net.loss(X, y, reg=0.1)[0]
    param_grad_num = eval_numerical_gradient(f, net.params[param_name], verbose=False)
    print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, grads[param_name])))

W2 max relative error: 8.048892e-10
b2 max relative error: 5.553999e-11
W1 max relative error: 1.121575e-08
b1 max relative error: 2.035406e-06

In [ ]:
```

[5pts] Q2.3 Train the Sigmoid network

To train the network we will use stochastic gradient descent (SGD), implemented in `TwoLayerNet.train`. Then we train a two-layer network on toy data.

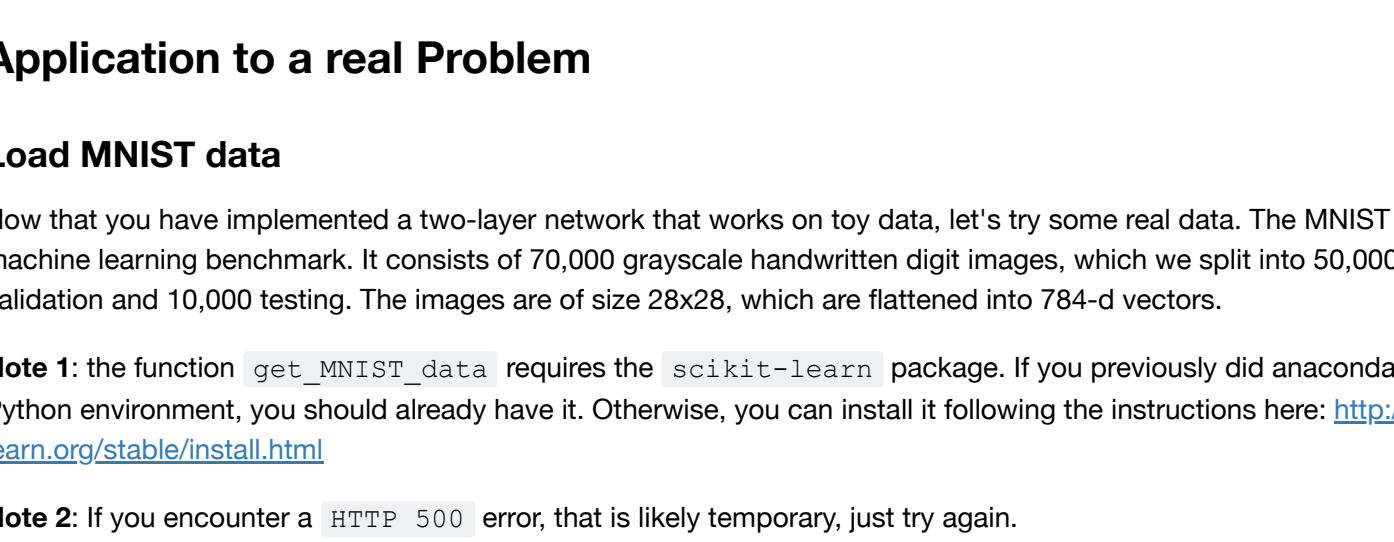
- Implement the prediction function `TwoLayerNet.predict`, which is called during training to keep track of training and validation accuracy.

You should get the first training loss around 0.1, which is good, but not too great for such a toy problem. One problem is that the gradient magnitude for `W1` (the first layer weights) stays small all the time, and the neural net doesn't get much "learning signals". This has to do with the saturation problem of the sigmoid activation function.

```
In [5]: net = init_toy_model('sigmoid', std=1e-1)
states = net.train(X, y, X_val, y_val,
                  learning_rate=0.5, reg=1e-5,
                  num_epochs=100, verbose=False)
print('Final training loss: ', states['loss_history'][-1])

# plot the loss history and gradient magnitudes
fig, (ax1, ax2) = plt.subplots(2, 1)
ax1.plot(states['loss_history'])
ax1.set_xlabel('epoch')
ax1.set_ylabel('training loss')
ax1.set_title('Training loss history')
ax2.plot(states['grad_magnitude_history'])
ax2.set_xlabel('iteration')
ax2.set_ylabel('||\nabla_{W1}||_1')
ax2.set_title('Gradient magnitude history ' + r'\nabla_{W1}_1')
ax2.set_ylim(0,1)
fig.tight_layout()
plt.show()

Final training loss: 0.10926794610680679
```



Application to a real Problem

Load MNIST data

Now that you have implemented a two-layer network that works on toy data, let's try some real data. The MNIST dataset is a standard machine learning benchmark. It consists of 70,000 grayscale handwritten digit images, which we split into 50,000 training, 10,000 validation and 10,000 testing. The images are of size 28x28, which are flattened into 784-d vectors.

Note 1: the function `get_MNIST_data` requires the `scikit-learn` package. If you previously did not install it, you should set up your Python environment, you should already have it. Otherwise, you can install it following the instructions here: <http://scikit-learn.org/stable/install.html>

Note 2: If you encounter a `HTTP 500 error`, that is likely temporary, just try again.

Note 3: Ensure that the downloaded MNIST file is 55.4MB (smaller file-sizes could indicate an incomplete download - which is possible)

```
In [6]: # Load MNIST
from util import get_MNIST_data

X_train, y_train, X_val, y_val, X_test, y_test = get_MNIST_data()
print('Train data shape: ', y_val, y_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

Train data shape: (50000, 784)
Train labels shape: (50000,)
Validation data shape: (10000, 784)
Validation labels shape: (10000,)
Test data shape: (10000, 784)
Test labels shape: (10000,)
```

Train a network on MNIST

We will now train a network on MNIST with 64 hidden units in the hidden layer. We train it using SGD, and decrease the learning rate with an exponential rate over time. This is achieved by multiplying the learning rate with a constant factor `learning_rate_decay` (which is less than 1) after each epoch. In effect, we are using a high learning rate initially, which is good for exploring the solution space, and using low learning rates later to encourage convergence to a local minimum (or [saddle point](#), which may happen more often).

- Train your MNIST network with 2 different activation functions: sigmoid and ReLU.

We first define some variables and utility functions. The `plot_states` function plots the histories of gradient magnitude, training loss, and accuracies on the training and validation sets. The `show_net_weights` function visualizes the weights learned in the first layer of the network. In most neural networks trained on visual data, the first layer weights typically show some visible structure when visualized. Both functions help you to diagnose the training process.

```
In [9]: input_size = 28 * 28
hidden_size = 64
num_classes = 10

# Plot the loss, accuracy and training / validation accuracies
def plot_states(states):
    fig, (ax1, ax2, ax3) = plt.subplots(3, 1)
    ax1.plot(states['grad_magnitude_history'])
    ax1.set_xlabel('iteration')
    ax1.set_ylabel('||\nabla_{W1}||_1')
    ax1.set_title('Gradient magnitude history ' + r'\nabla_{W1}_1')
    ax1.set_ylim(0, np.minimum(100, np.max(states['grad_magnitude_history'])))
    ax2.plot(states['loss_history'])
    ax2.set_xlabel('epoch')
    ax2.set_ylabel('loss')
    ax2.set_ylim(0, 100)
    ax3.plot(states['train_acc_history'], label='train')
    ax3.plot(states['val_acc_history'], label='val')
    ax3.set_xlabel('epoch')
    ax3.set_ylabel('Classification accuracy')
    fig.tight_layout()
    plt.show()

# Visualize the weights of the network
def show_net_weights(grid):
    W1 = net.params['W1']
    W1 = W1.T.reshape(-1, 28, 28)
    plt.imshow(show_weights_grid(W1, padding=3).astype('uint8'))
    plt.gca().axis('off')
    plt.show()
```

Sigmoid network

```
In [10]: sigmoid_net = TwoLayerMLP(input_size, hidden_size, num_classes, activation='sigmoid', std=1e-1)

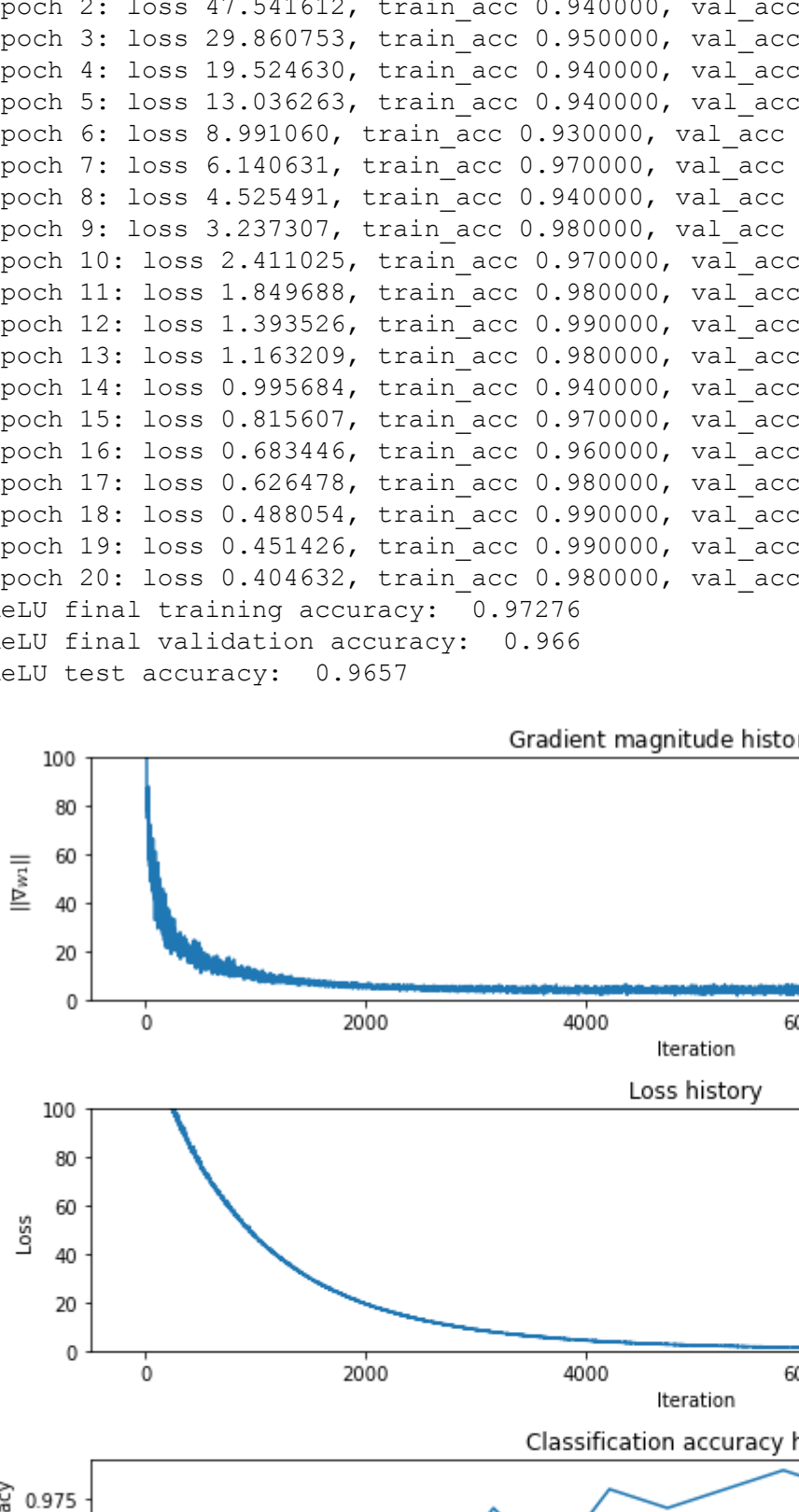
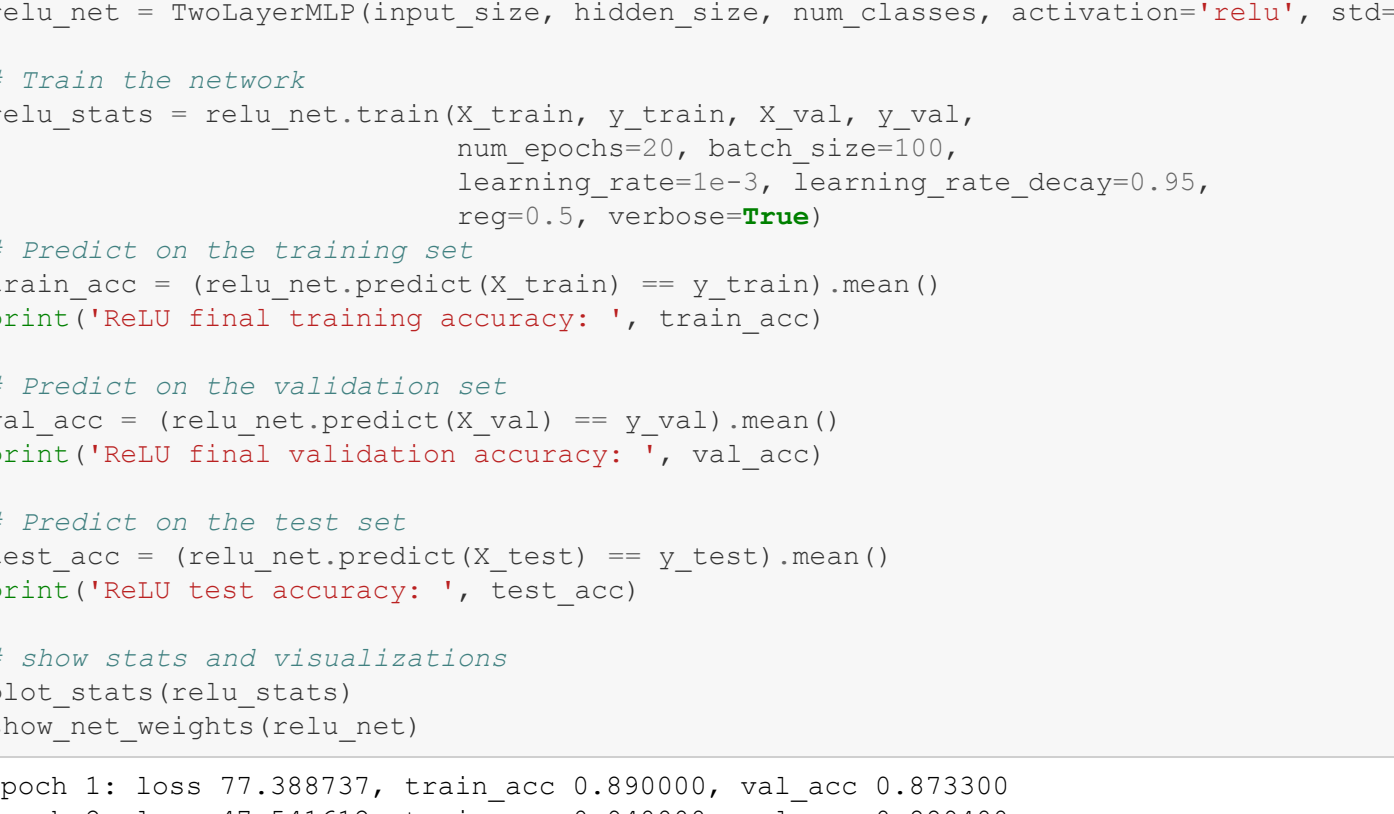
# Train the network
sigmoid_state = sigmoid_net.train(X_train, y_train, X_val, y_val,
                                 num_epochs=20, batch_size=100,
                                 learning_rate=0.5, learning_rate_decay=0.95,
                                 reg=0.5, verbose=True)

# Predict on the training set
train_acc = (sigmoid_net.predict(X_train) == y_train).mean()
print('Sigmoid final training accuracy: ', train_acc)

# Predict on the validation set
val_acc = (sigmoid_net.predict(X_val) == y_val).mean()
print('Sigmoid final validation accuracy: ', val_acc)

# Predict on the test set
test_acc = (sigmoid_net.predict(X_test) == y_test).mean()
print('Sigmoid test accuracy: ', test_acc)

# show stats and visualizations
plot_states(sigmoid_net)
show_net_weights(sigmoid_net)
```



ReLU network

```
In [11]: relu_net = TwoLayerMLP(input_size, hidden_size, num_classes, activation='relu', std=1e-1)

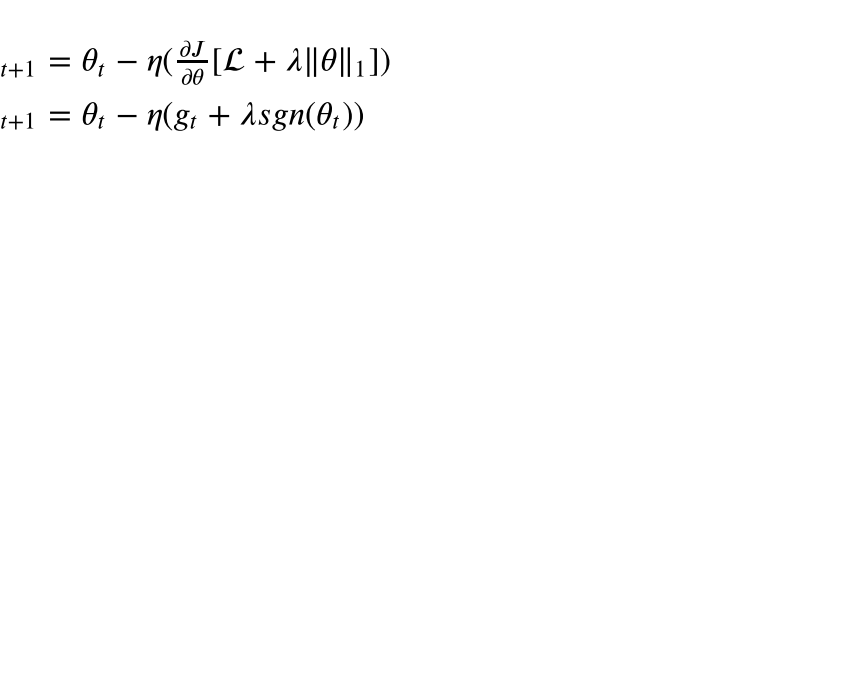
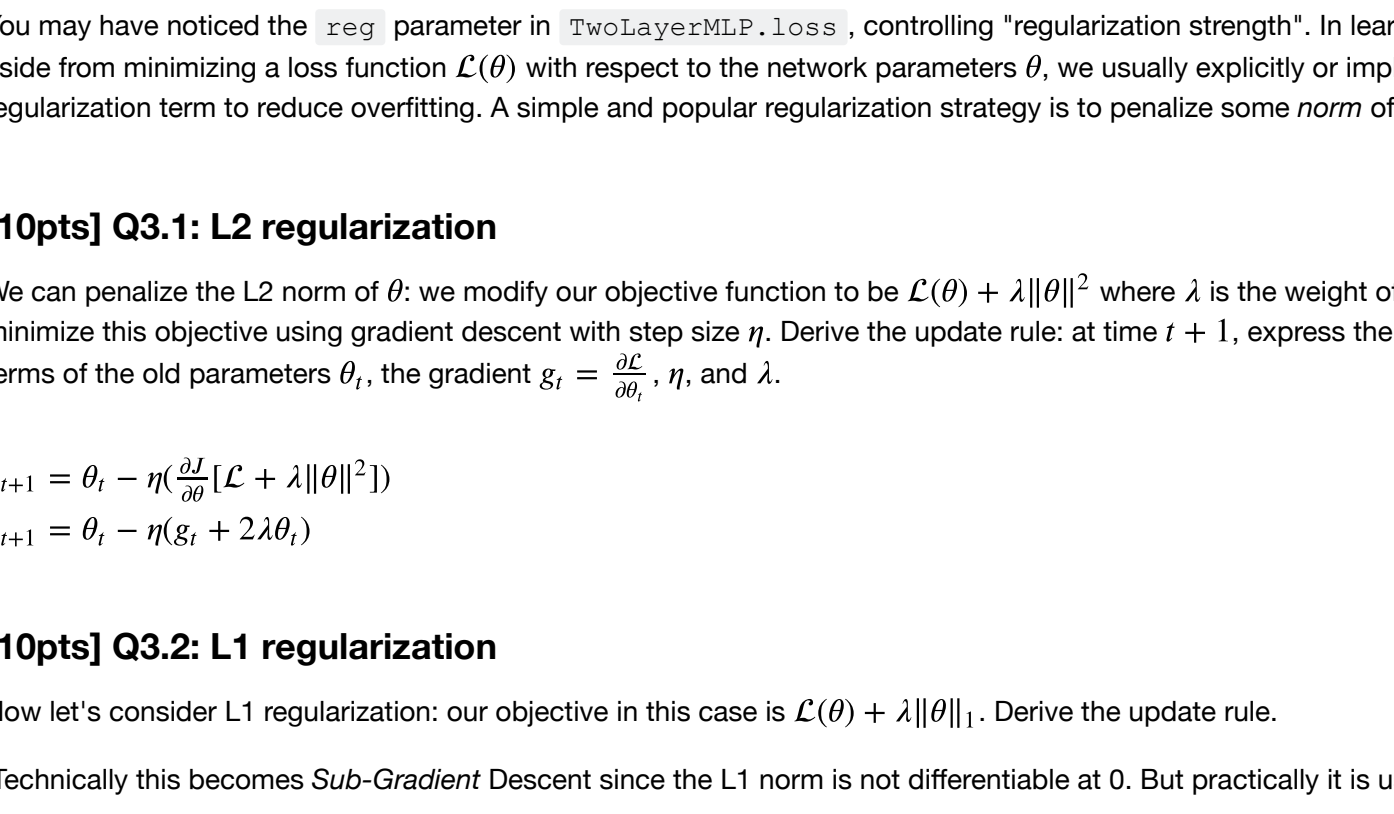
# Train the network
relu_state = relu_net.train(X_train, y_train, X_val, y_val,
                           num_epochs=20, batch_size=100,
                           learning_rate=0.5, learning_rate_decay=0.95,
                           reg=0.5, verbose=True)

# Predict on the training set
train_acc = (relu_net.predict(X_train) == y_train).mean()
print('ReLU final training accuracy: ', train_acc)

# Predict on the validation set
val_acc = (relu_net.predict(X_val) == y_val).mean()
print('ReLU final validation accuracy: ', val_acc)

# Predict on the test set
test_acc = (relu_net.predict(X_test) == y_test).mean()
print('ReLU test accuracy: ', test_acc)

# show stats and visualizations
plot_states(relu_state)
show_net_weights(relu_net)
```



[5pts] Q2.5

Which activation function would you choose in practice? Why?

ReLU because it provides a much higher accuracy percentage.

[20pts] Problem 3: Simple Regularization Methods

You may have noticed the `reg` parameter in `TwoLayerMLP.loss`, controlling "regularization strength". In learning neural networks, aside from minimizing the loss function $\mathcal{L}(\theta)$ with respect to the network parameters θ , we usually explicitly or implicitly add some regularization term to reduce overfitting. A simple and popular regularization strategy is to penalize some *norm* of θ .

[10pts] Q3.1: L2 regularization

We can penalize the L2 norm of θ : we modify our objective function to be $\mathcal{L}(\theta) + \lambda \|\theta\|^2$ where λ is the weight of regularization. We will minimize this objective using gradient descent with step size η . Derive the update rule: at time $t + 1$, express the new parameters θ_{t+1} in terms of the old parameters θ_t , the gradient $g_t = \frac{\partial}{\partial \theta} \mathcal{L}$, and λ .

$$\theta_{t+1} = \theta_t - \eta \left(\frac{\partial}{\partial \theta} [\mathcal{L} + \lambda \|\theta\|^2] \right)$$
$$\theta_{t+1} = \theta_t - \eta (g_t + 2\lambda \theta_t)$$

[10pts] Q3.2: L1 regularization

Now let's consider L1 regularization: our objective in this case is $\mathcal{L}(\theta) + \lambda \|\theta\|_1$. Derive the update rule.

(Technically this becomes *Sub-Gradient* Descent since the L1 norm is not differentiable at 0. But practically it is usually not an issue.)

$$\theta_{t+1} = \theta_t - \eta \frac{\partial}{\partial \theta} [\mathcal{L} + \lambda \|\theta\|_1]$$
$$\theta_{t+1} = \theta_t - \eta (g_t + \lambda \text{sgn}(\theta_t))$$