

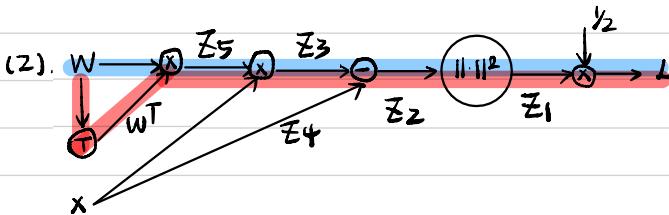
1. (15 points) **Backpropagation for autoencoders.** In an autoencoder, we seek to reconstruct the original data after some operation that reduces the data's dimensionality. We may be interested in reducing the data's dimensionality to gain a more compact representation of the data.

For example, consider $\mathbf{x} \in \mathbb{R}^n$. Further, consider $\mathbf{W} \in \mathbb{R}^{m \times n}$ where $m < n$. Then \mathbf{Wx} is of lower dimensionality than \mathbf{x} . One way to design \mathbf{W} so that \mathbf{Wx} still contains key features of \mathbf{x} is to minimize the following expression

$$\mathcal{L} = \frac{1}{2} \|\mathbf{W}^T \mathbf{Wx} - \mathbf{x}\|^2$$

with respect to \mathbf{W} . (To be complete, autoencoders also have a nonlinearity in each layer, i.e., the loss is $\frac{1}{2} \|f(\mathbf{W}^T f(\mathbf{Wx})) - \mathbf{x}\|^2$. However, we'll work with the linear example.)

(1) Because if $L \rightarrow 0$, then it means $\mathbf{W}^T \mathbf{Wx} - \mathbf{x} \rightarrow 0$, regard $\mathbf{W}^T \mathbf{W}$ as a function $f(\cdot)$. Then $f(\mathbf{x}) \rightarrow \mathbf{x}$, then it means after transformation, $f(\mathbf{x})$ is still very close to original \mathbf{x} , it means information of \mathbf{x} has been preserved.



(2). There are \mathbf{w} and \mathbf{w}^T two paths to account for $\frac{\partial L}{\partial \mathbf{w}}$.

$$\frac{\partial L}{\partial \mathbf{w}} = (\frac{\partial L}{\partial \mathbf{w}})_{\text{blue}} + (\frac{\partial L}{\partial \mathbf{w}^T})_{\text{pink}}^T$$

$$(3). L = \frac{1}{2} z_1, \frac{\partial L}{\partial z_1} = \frac{1}{2}, z_1 = \|z_2\|^2, \frac{\partial z_1}{\partial z_2} = 2z_2, \frac{\partial L}{\partial z_2} = \frac{\partial z_1}{\partial z_2} \cdot \frac{\partial L}{\partial z_1} = z_2$$

"Add" passes gradient: $\frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial z_2}$, $z_3 = z_5 \cdot x$, $\frac{\partial L}{\partial z_5} = \frac{\partial z_3}{\partial z_5} \cdot \frac{\partial L}{\partial z_3} = \frac{\partial L}{\partial z_3} \cdot x^T$

$$(\frac{\partial L}{\partial \mathbf{w}})_1 = \frac{\partial z_5}{\partial \mathbf{w}} \cdot \frac{\partial L}{\partial z_5} = \mathbf{w} \cdot \frac{\partial L}{\partial z_5} = \mathbf{w} \cdot z_2 \cdot x^T, (\frac{\partial L}{\partial \mathbf{w}^T})_2 = \frac{\partial z_5}{\partial \mathbf{w}^T} \frac{\partial L}{\partial z_5} = \frac{\partial L}{\partial z_5} \cdot \mathbf{w}^T = z_2 \cdot x^T \cdot \mathbf{w}^T$$

$$z_5 = \mathbf{w}^T \mathbf{w}$$

$$\Rightarrow \frac{\partial L}{\partial \mathbf{w}} = (\frac{\partial L}{\partial \mathbf{w}})_1 + (\frac{\partial L}{\partial \mathbf{w}^T})_2^T = \mathbf{w} \cdot z_2 \cdot x^T + \mathbf{w} \cdot x \cdot z_2^T = \mathbf{w} \left((\mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x}) \mathbf{x}^T + \mathbf{x} (\mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x})^T \right)$$

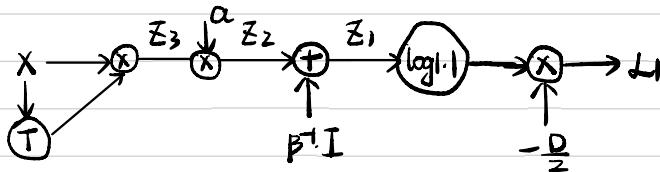
$$z_2 = \mathbf{w}^T \mathbf{w} \mathbf{x} - \mathbf{x}$$

2. $\mathcal{L} = -c - \frac{D}{2} \log |\mathbf{K}| - \frac{1}{2} \text{tr}(\mathbf{K}^{-1} \mathbf{Y} \mathbf{Y}^T)$ $\mathbf{K} = \alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}$

$$\mathcal{L}_1 = -\frac{D}{2} \log |\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I}|$$

$$\mathcal{L}_2 = -\frac{1}{2} \text{tr}((\alpha \mathbf{X} \mathbf{X}^T + \beta^{-1} \mathbf{I})^{-1} \mathbf{Y} \mathbf{Y}^T)$$

(1). Draw Computation Graph for \mathcal{L}_1 .



(2) Compute $\frac{\partial \mathcal{L}_1}{\partial \mathbf{x}}$: $\mathcal{L}_1 = -\frac{D}{2} \log |\det \mathbf{Z}_1|$, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1} = -\frac{D}{2} (\mathbf{Z}_1^T)^{-1} = -\frac{D}{2} (\mathbf{Z}_1^T)^{-1}$.
 $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_2} = \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1}$, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_3} = \alpha \cdot \frac{\partial \mathcal{L}_1}{\partial \mathbf{z}_1} = -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1}$

$$\mathbf{Z}_3 = \mathbf{X} \mathbf{X}^T, \frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = \frac{\partial \mathbf{Z}_3}{\partial \mathbf{X}}, \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} = \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} \cdot (\mathbf{X}^T)^T = \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} \cdot \mathbf{X}$$

$$\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}^T} = \frac{\partial \mathbf{Z}_3}{\partial \mathbf{X}^T}, \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} = \mathbf{X}^T \cdot \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3}$$

$$\Rightarrow \frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} \right) + \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}^T} \right)^T = \frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} \cdot \mathbf{X} + \left(\frac{\partial \mathcal{L}_1}{\partial \mathbf{Z}_3} \right)^T \cdot \mathbf{X}$$

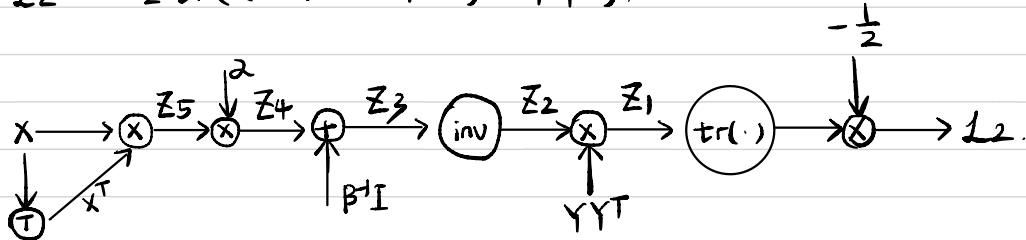
$$= -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1} \cdot \mathbf{X} + \left(-\frac{D\alpha}{2} \cdot (\mathbf{Z}_1^T)^T \right)^T \cdot \mathbf{X}$$

$\mathbf{Z}_1 = \mathbf{X} \mathbf{X}^T \alpha + \beta^{-1} \mathbf{I}$ is symmetric, $\frac{\partial \mathcal{L}_1}{\partial \mathbf{X}} = -\frac{D\alpha}{2} (\mathbf{Z}_1^T)^{-1} \mathbf{X} - \frac{D\alpha}{2} \cdot \mathbf{Z}_1^{-1} \mathbf{X} = -D\alpha \mathbf{Z}_1^{-1} \mathbf{X}$

$$= -D\alpha \cdot (\mathbf{X} \mathbf{X}^T \alpha + \beta^{-1} \mathbf{I})^{-1} \mathbf{X}$$

(3). Draw Computational Graph for L_2 .

$$L_2 = -\frac{1}{2} \text{tr}((\alpha XX^T + \beta^{-1} I)^{-1} \cdot YY^T).$$



(4) Compute $\frac{\partial L_2}{\partial X}$.

$$L_2 = -\frac{1}{2} \text{tr}(Z_1), \quad \frac{\partial L_2}{\partial Z_1} = -\frac{1}{2} \frac{(Z_1)}{\partial Z_1} = -\frac{1}{2} I$$

$$Z_1 = Z_2(YY^T), \quad \frac{\partial L_2}{\partial Z_2} = \frac{\partial Z_2}{\partial Z_2} \cdot \frac{\partial L_2}{\partial Z_1} = -\frac{1}{2} I \cdot (YY^T)^T = -\frac{1}{2} YY^T.$$

$$Z_2 = (Z_3)^{-1}, \quad \frac{\partial Z_2}{\partial Z_3} = -Z_3^{-T} \cdot \frac{\partial Z_2}{\partial Z_3^{-1}} \cdot Z_3^{-T} = -Z_3^{-T} \cdot I \cdot Z_3^{-T} = -(Z_3^{-T})^2.$$

$$\frac{\partial L}{\partial Z_3} = \frac{\partial Z_2}{\partial Z_3} \cdot \frac{\partial L}{\partial Z_2} = + (Z_3^{-T})^2 \cdot \frac{1}{2} YY^T$$

"Add" passes gradient: $\frac{\partial L_2}{\partial Z_4} = \frac{\partial L_2}{\partial Z_3}$

$$Z_4 = \alpha \cdot Z_5, \quad \frac{\partial L_2}{\partial Z_5} = \alpha \cdot \frac{\partial L_2}{\partial Z_4} = \alpha \cdot \frac{\partial L_2}{\partial Z_3} = \frac{\alpha}{2} \cdot (Z_3^{-T} \cdot Z_3^{-T}) \cdot YY^T$$

$$Z_5 = XX^T, \quad \frac{\partial L_2}{\partial X} = \frac{\partial Z_5}{\partial X} \cdot \frac{\partial L_2}{\partial Z_5} = \frac{\partial L_2}{\partial Z_5} \cdot (X^T)^T = \frac{\partial L_2}{\partial Z_5} \cdot X$$

$$\frac{\partial L_2}{\partial X^T} = \frac{\partial Z_5}{\partial X^T} \cdot \frac{\partial L_2}{\partial Z_5} = X^T \cdot \frac{\partial L_2}{\partial Z_5}$$

$$\frac{\partial L_2}{\partial X} = (\frac{\partial L_2}{\partial X})_1 + (\frac{\partial L_2}{\partial X^T})^T = \frac{\partial L_2}{\partial Z_5} \cdot X + (\frac{\partial L_2}{\partial Z_5})^T \cdot X$$

$$= \frac{\alpha}{2} \cdot (Z_3^{-T} \cdot Z_3^{-T}) \cdot YY^T \cdot X + \left(\frac{\alpha}{2} \cdot (Z_3^{-T})^2 \cdot YY^T \right)^T \cdot X$$

$$= \frac{\alpha}{2} \left((Z_3^{-T})^2 YY^T + YY^T \cdot (Z_3^{-T})^2 \right) X$$

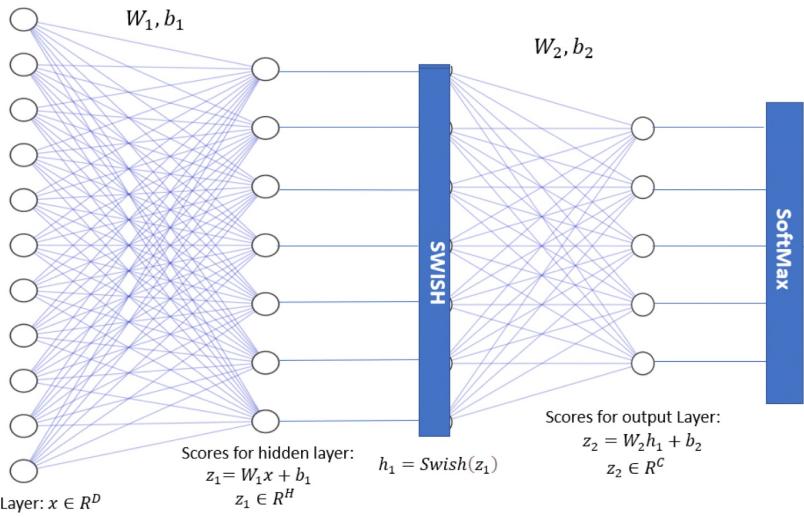
where $Z_3 = \alpha XX^T + \beta^{-1} I$

$$(e). \frac{\partial L}{\partial X} = \frac{\partial}{\partial X} (-C + L_1 + L_2) = \frac{\partial L_1}{\partial X} + \frac{\partial L_2}{\partial X}$$

$$= -D \cdot \alpha \cdot (XX^T \cdot \alpha + \beta^{-1} I)^{-1} X$$

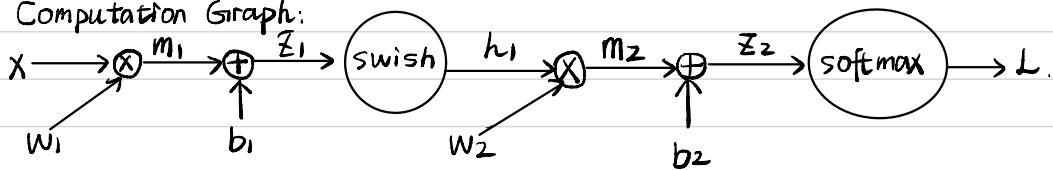
$$+ \frac{\alpha}{2} \left((Z_3^{-T})^2 YY^T + YY^T \cdot (Z_3^{-T})^2 \right) X$$

$$Z_3 = \alpha XX^T + \beta^{-1} I.$$



$$\text{Swish}(k) = k \cdot \sigma(k)$$

(a). Computation Graph:



(b). m_1, m_2 defined as above.

“Add” passes gradient: $\frac{\partial L}{\partial m_2} = \frac{\partial L}{\partial z_2}, \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2}$.

$$m_2 = W_2 h_1, \frac{\partial L}{\partial W_2} = \frac{\partial m_2}{\partial W_2}, \frac{\partial L}{\partial m_2} = \frac{\partial L}{\partial m_2} \cdot h_1^T = \frac{\partial L}{\partial z_2} \cdot h_1^T$$

$$(c) \frac{\partial L}{\partial h_1} = \frac{\partial m_2}{\partial h_1} \frac{\partial L}{\partial m_2} = W_2^T \frac{\partial L}{\partial z_2}$$

$$h_1 = \text{swish}(z_1), \frac{\partial h_1}{\partial z_1} = \sigma(z_1) + z_1(1 - \sigma(z_1))$$

$$= z_1 \cdot \sigma(z_1)$$

$$\frac{\partial L}{\partial z_1} = \frac{\partial h_1}{\partial z_1} \cdot \frac{\partial L}{\partial h_1} = \text{diag}(\sigma(z_1) + z_1(1 - \sigma(z_1))) \frac{\partial L}{\partial h_1} = \text{diag}(\sigma(z_1) + z_1(1 - \sigma(z_1))) \cdot W_2^T \frac{\partial L}{\partial z_2}$$

$$\text{“Add” passes gradient: } \frac{\partial}{\partial b_1} = \frac{\partial L}{\partial z_1}, \frac{\partial L}{\partial m_1} = \frac{\partial L}{\partial z_1}$$

$$m_1 = W_1 x, \frac{\partial L}{\partial W_1} = \frac{\partial m_1}{\partial W_1} \frac{\partial L}{\partial m_1} = \frac{\partial L}{\partial m_1} \cdot x^T = \frac{\partial L}{\partial z_1} \cdot x^T$$

two_layer_nn

February 2, 2023

0.1 This is the 2-layer neural network notebook for ECE C147/C247 Homework #3

Please follow the notebook linearly to implement a two layer neural network.

Please print out the notebook entirely when completed.

The goal of this notebook is to give you experience with training a two layer neural network.

```
[ ]: from google.colab import drive  
drive.mount('/content/drive')
```

Mounted at /content/drive

```
[ ]: import sys  
sys.path.append('/content/drive/MyDrive/247-NeuralNetwork/247-hw3/HW3_code')  
import random  
import numpy as np  
from utils.data_utils import load_CIFAR10  
import matplotlib.pyplot as plt  
  
%matplotlib inline  
%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))))
```

0.2 Toy example

Before loading CIFAR-10, there will be a toy example to test your implementation of the forward and backward pass. Make sure to read the description of TwoLayerNet class in neural_net.py file , understand the architecture and initializations

```
[ ]: from nnndl.neural_net import TwoLayerNet
```

```
[ ]: # 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():
    np.random.seed(0)
    return TwoLayerNet(input_size, hidden_size, num_classes, std=1e-1)

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

net = init_toy_model()
X, y = init_toy_data()

```

[]: `print(X,y)`

```

[[ 16.24345364 -6.11756414 -5.28171752 -10.72968622]
 [  8.65407629 -23.01538697  17.44811764 -7.61206901]
 [  3.19039096 -2.49370375  14.62107937 -20.60140709]
 [ -3.22417204 -3.84054355  11.33769442 -10.99891267]
 [ -1.72428208 -8.77858418   0.42213747   5.82815214]] [0 1 2 2 1]

```

0.2.1 Compute forward pass scores

[]: `## Implement the forward pass of the neural network.`
`## See the loss() method in TwoLayerNet class for the same`

```

# Note, there is a statement if y is None: return scores, which is why
# the following call will calculate the scores.
scores = net.loss(X)
print('Your scores:')
print(scores)
print()
print('correct scores:')
correct_scores = np.asarray([
    [-1.07260209,  0.05083871, -0.87253915],
    [-2.02778743, -0.10832494, -1.52641362],
    [-0.74225908,  0.15259725, -0.39578548],
    [-0.38172726,  0.10835902, -0.17328274],
    [-0.64417314, -0.18886813, -0.41106892]])
print(correct_scores)
print()

```

```
# The difference should be very small. We get < 1e-7
print('Difference between your scores and correct scores:')
print(np.sum(np.abs(scores - correct_scores)))
```

Your scores:

```
[[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]
```

correct scores:

```
[[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]
```

Difference between your scores and correct scores:
3.3812311957259755e-08

0.2.2 Forward pass loss

```
[ ]: loss, _ = net.loss(X, y, reg=0.05)
correct_loss = 1.071696123862817

# should be very small, we get < 1e-12
print("Loss:", loss)
print('Difference between your loss and correct loss:')
print(np.sum(np.abs(loss - correct_loss)))
```

Loss: 1.071696123862817
Difference between your loss and correct loss:
0.0

0.2.3 Backward pass

Implements the backwards pass of the neural network. Check your gradients with the gradient check utilities provided.

```
[ ]: from utils.gradient_check import eval_numerical_gradient

# 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.

loss, grads = net.loss(X, y, reg=0.05)
```

```
# these should all be less than 1e-8 or so
for param_name in grads:
    f = lambda W: net.loss(X, y, reg=0.05)[0]
    param_grad_num = eval_numerical_gradient(f, net.params[param_name], □
    ↪verbose=False)
    print('{} max relative error: {}'.format(param_name, □
    ↪rel_error(param_grad_num, grads[param_name])))
```

```
W2 max relative error: 2.963221034761121e-10
b2 max relative error: 1.839165909006465e-10
b1 max relative error: 3.1726806716844575e-09
W1 max relative error: 1.283285096965795e-09
```

0.2.4 Training the network

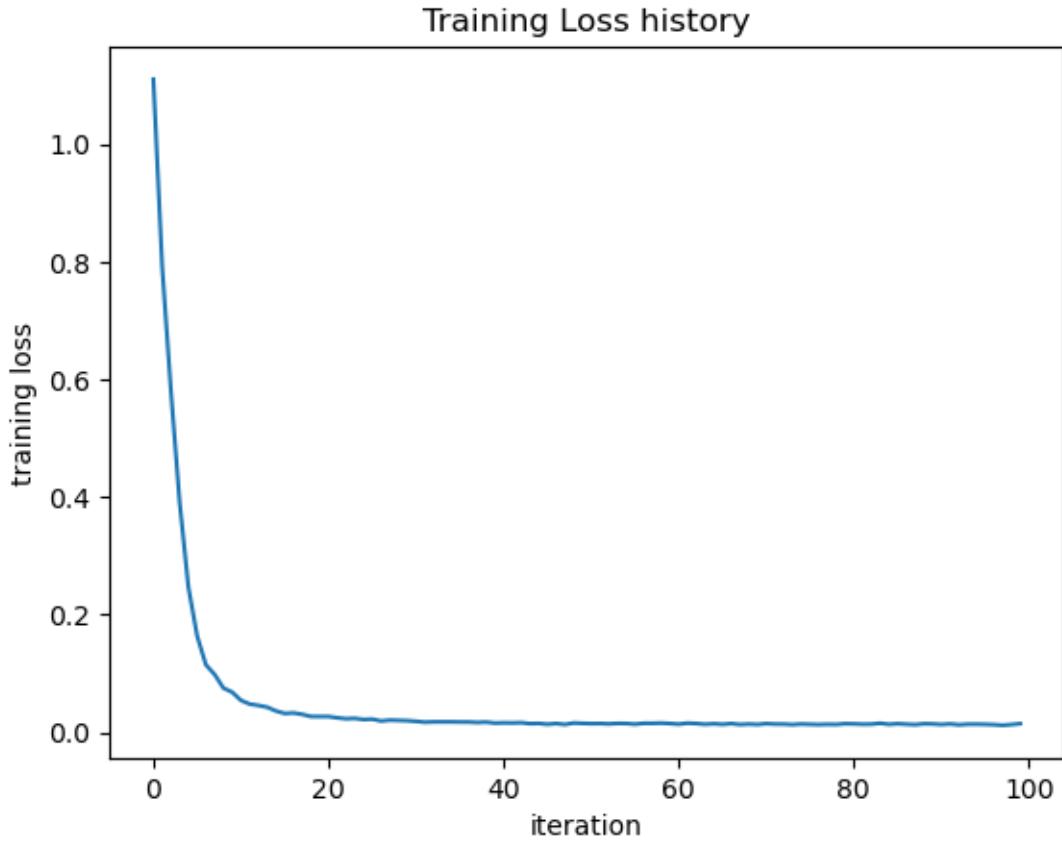
Implement `neural_net.train()` to train the network via stochastic gradient descent, much like the softmax and SVM.

```
[ ]: net = init_toy_model()
stats = net.train(X, y, X, y,
                  learning_rate=1e-1, reg=5e-6,
                  num_iters=100, verbose=False)

print('Final training loss: ', stats['loss_history'][-1])

# plot the loss history
plt.plot(stats['loss_history'])
plt.xlabel('iteration')
plt.ylabel('training loss')
plt.title('Training Loss history')
plt.show()
```

```
Final training loss: 0.014497864587765875
```



0.3 Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

```
[ ]: from utils.data_utils import load_CIFAR10

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net classifier.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = '/Users/yiwenzhang/Desktop/23winter/247 neural network/hw2/
    ↵cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # Subsample the data
    mask = list(range(num_training, num_training + num_validation))
    X_val = X_train[mask]
    y_val = y_train[mask]
```

```

mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]

# Normalize the data: subtract the mean image
mean_image = np.mean(X_train, axis=0)
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image

# Reshape data to rows
X_train = X_train.reshape(num_training, -1)
X_val = X_val.reshape(num_validation, -1)
X_test = X_test.reshape(num_test, -1)

return X_train, y_train, X_val, y_val, X_test, y_test

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
print('Train data shape: ', X_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: (49000, 3072)
Train labels shape: (49000,)
Validation data shape: (1000, 3072)
Validation labels shape: (1000,)
Test data shape: (1000, 3072)
Test labels shape: (1000,)

```

0.3.1 Running SGD

If your implementation is correct, you should see a validation accuracy of around 28-29%.

```

[ ]: input_size = 32 * 32 * 3
hidden_size = 50
num_classes = 10
net = TwoLayerNet(input_size, hidden_size, num_classes)

# Train the network
stats = net.train(X_train, y_train, X_val, y_val,

```

```

    num_iters=1000, batch_size=200,
    learning_rate=1e-4, learning_rate_decay=0.95,
    reg=0.25, verbose=True)

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

# Save this net as the variable subopt_net for later comparison.
subopt_net = net

```

```

iteration 0 / 1000: loss 2.302757518613176
iteration 100 / 1000: loss 2.3021201657167145
iteration 200 / 1000: loss 2.2956136292668456
iteration 300 / 1000: loss 2.2518260278684843
iteration 400 / 1000: loss 2.188995574979717
iteration 500 / 1000: loss 2.1162531781593175
iteration 600 / 1000: loss 2.0646712637358062
iteration 700 / 1000: loss 1.990168856378208
iteration 800 / 1000: loss 2.0028266168788966
iteration 900 / 1000: loss 1.946516508873091
Validation accuracy: 0.283

```

0.4 Questions:

The training accuracy isn't great.

- (1) What are some of the reasons why this is the case? Take the following cell to do some analyses and then report your answers in the cell following the one below.
- (2) How should you fix the problems you identified in (1)?

```
[ ]: stats['train_acc_history']
```

```
[ ]: [0.095, 0.15, 0.25, 0.25, 0.315]
```

```

[ ]: from pickle import stat
# ===== #
# YOUR CODE HERE:
#   Do some debugging to gain some insight into why the optimization
#   isn't great.
# ===== #

# Plot the loss function and train / validation accuracies
plt.plot(stats['loss_history'])
plt.title('Loss History over Iteration')
plt.xlabel('Iter')
plt.ylabel('Loss')
plt.show()

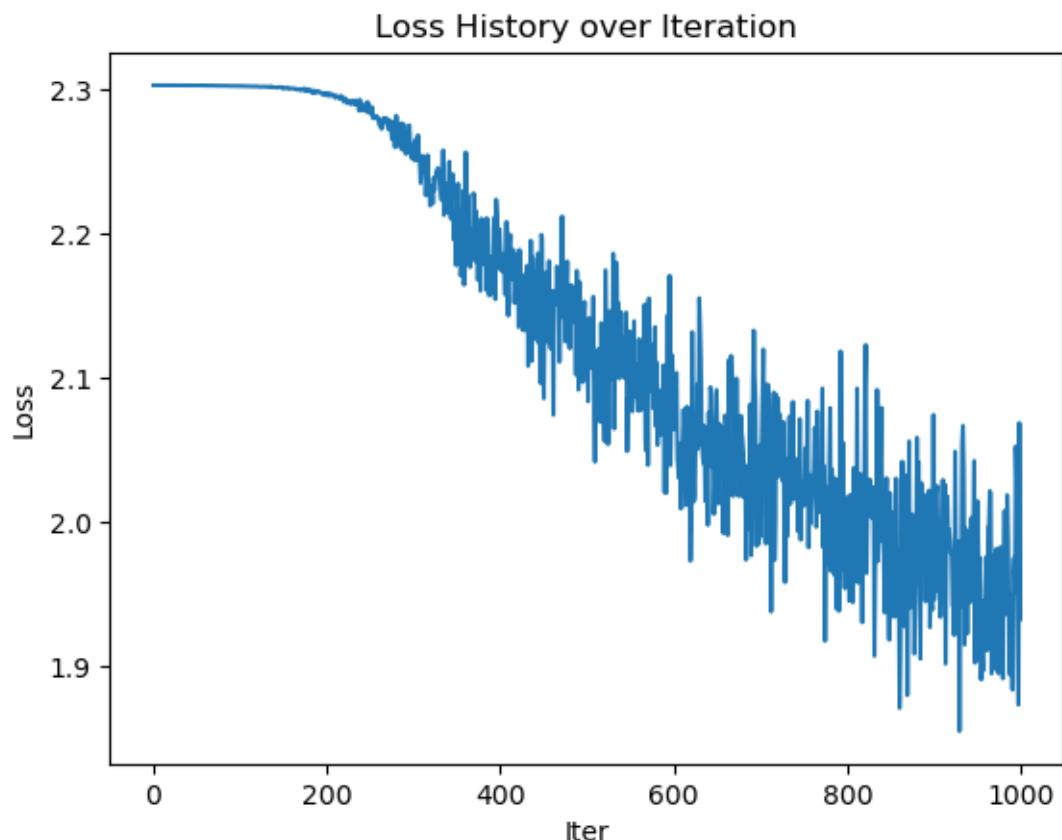
```

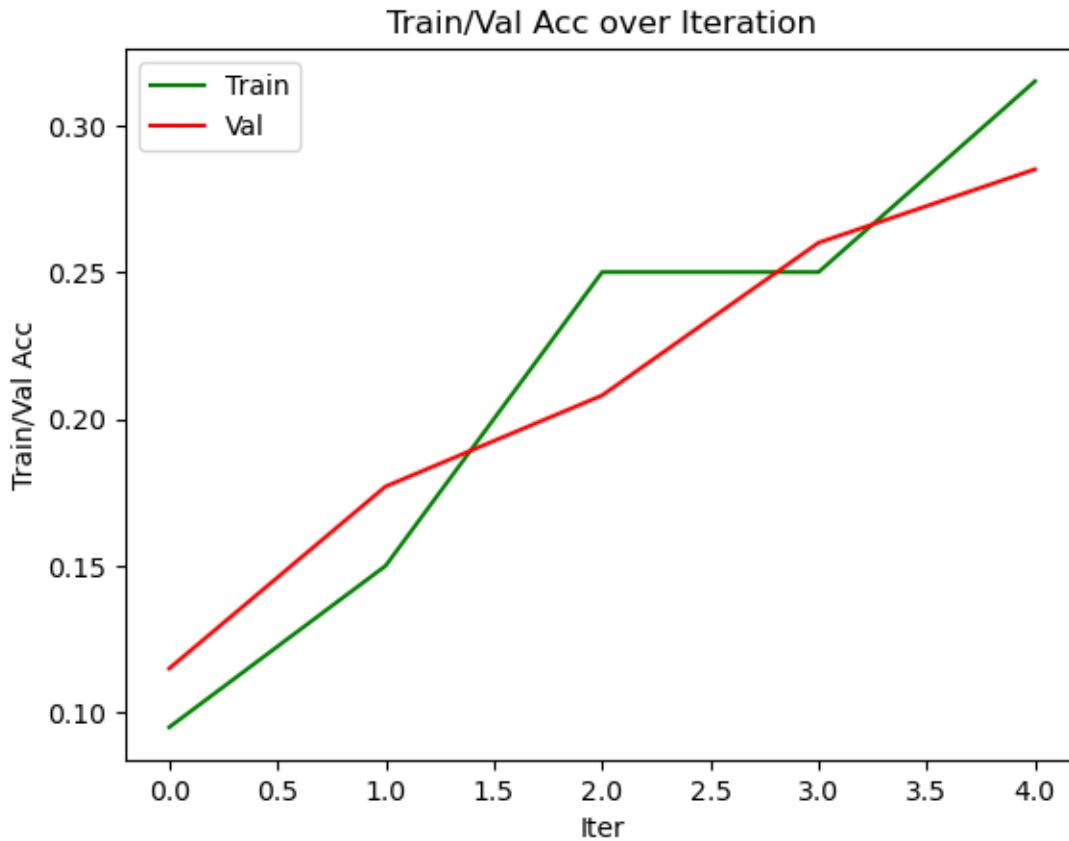
```

plt.plot(stats['train_acc_history'], label = 'Train', color = 'green')
plt.plot(stats['val_acc_history'], label = 'Val',color = 'red')
plt.title('Train/Val Acc over Iteration')
plt.xlabel('Iter')
plt.ylabel('Train/Val Acc')
plt.legend(loc=0)
plt.show()

pass
# ===== #
# END YOUR CODE HERE
# ===== #

```





0.5 Answers:

- (1) The training accuracy shows an increasing trend at the end of iteration, it means our learning rate is too small which causes our model takes so many interations but still doesn't converge.

Another reason could be that our model is too simple (like only two layers) to do the image classification.

The reason could also be that the regularization is too strong that it's difficult for our model to learn.

- (2) We could try to increase the step size of learning rate and try to add more layers to make the model stronger and more capacity so that it can handle better. Or we could decrease the regularization (not decrease too much otherwise it might cause overfitting) And we could also increase the num_iterations and enlarge the batch_size so that our net could learn longer.

0.6 Optimize the neural network

Use the following part of the Jupyter notebook to optimize your hyperparameters on the validation set. Store your nets as best_net.

Parameters which been optimizaed:

learning rate : 0.0001, 0.0005, 0.001
regularization: 0.05, 0.1, 0.2, 0.3, 0.5
batch_size: 200, 300, 400, 500
num_interation: 2000

Because the optimization of num_iteration here only have 2 options and it usually takes a long time, I split the whole optimization into 2 parts: the 1st is processed by num_ite = 1000, the 2scd is processed by num_ite = 2000.

```
[ ]: # When num_iteration is 2000
best_net = None # store the best model into this
# ===== #
# YOUR CODE HERE:
#   Optimize over your hyperparameters to arrive at the best neural
#   network. You should be able to get over 50% validation accuracy.
#   For this part of the notebook, we will give credit based on the
#   accuracy you get. Your score on this question will be multiplied by:
#       min(floor((X - 28%)) / %22, 1)
#   where if you get 50% or higher validation accuracy, you get full
#   points.
# Note, you need to use the same network structure (keep hidden_size = 50) !
# ===== #
import joblib
lr = [1e-4, 5*1e-4, 1e-3]
reg = [0.05, 0.1, 0.2, 0.3, 0.5]
batch_size = [200, 300, 400, 500]
val_acc_list = []
best_lr, best_reg, best_batch_size = 0, 0, 0
best_val_acc = 0
iter = 0
for i in range(len(lr)):
    for j in range(len(reg)):
        for k in range(len(batch_size)):
            net = TwoLayerNet(input_size, hidden_size, num_classes)
            best_model = net.train(X_train, y_train, X_val, y_val,
                num_iters=2000, batch_size=batch_size[k],
                learning_rate=lr[i], learning_rate_decay=0.95,
                reg=reg[j], verbose=True)
            val_acc = (net.predict(X_val) == y_val).mean()
            val_acc_list.append(val_acc)
            iter += 1
            print(iter, " Iteration")
            print('Learning rate: ', lr[i], "Reg: ", reg[j],
                  'Batch size: ', batch_size[k], 'Val_acc: ', val_acc)
            if val_acc > best_val_acc:
                best_val_acc = val_acc
```

```

        best_lr = lr[i]
        best_reg = reg[j]
        best_batch_size = batch_size[k]
        joblib.dump(net, 'best_model.dat')
        print( )
        print( )
best_net = joblib.load('best_model.dat')
print('Best Learning rate: ', best_lr)
print("Best Reg: ", best_reg)
print("Best Batch Size: ", best_batch_size)
print('Best Val_acc: ', best_val_acc)
pass

# ===== #
# END YOUR CODE HERE
# ===== #

val_acc = (best_net.predict(X_val) == y_val).mean()
if best_val_acc > 0.5:
    print('Validation accuracy: ', best_val_acc)
else:
    print("The best accuracy at num_iteration = 1000 is lower than 0.5")

```

```

iteration 0 / 2000: loss 2.302625426147543
iteration 100 / 2000: loss 2.302168765141486
iteration 200 / 2000: loss 2.2991812040710204
iteration 300 / 2000: loss 2.2641034954578934
iteration 400 / 2000: loss 2.217941694869426
iteration 500 / 2000: loss 2.1770821786296493
iteration 600 / 2000: loss 2.034967902437716
iteration 700 / 2000: loss 2.015380701398895
iteration 800 / 2000: loss 1.9730459987150968
iteration 900 / 2000: loss 2.0030313508883753
iteration 1000 / 2000: loss 1.8902312389513123
iteration 1100 / 2000: loss 1.8836828080267949
iteration 1200 / 2000: loss 2.0152404593949576
iteration 1300 / 2000: loss 1.8135211131982254
iteration 1400 / 2000: loss 1.8747462093197926
iteration 1500 / 2000: loss 1.8300951001302264
iteration 1600 / 2000: loss 1.8692129551574395
iteration 1700 / 2000: loss 1.8774043663892894
iteration 1800 / 2000: loss 1.859324923136817
iteration 1900 / 2000: loss 1.8034393504280757
1 Iteration
Learning rate: 0.0001 Reg: 0.05 Batch size: 200 Val_acc: 0.358

```

```
iteration 0 / 2000: loss 2.3026204826527223
```

```
iteration 100 / 2000: loss 1.9246812866173355
iteration 200 / 2000: loss 1.786948782697525
iteration 300 / 2000: loss 1.7088101968830318
iteration 400 / 2000: loss 1.6587468048406497
iteration 500 / 2000: loss 1.562652174021649
iteration 600 / 2000: loss 1.5904557907006343
iteration 700 / 2000: loss 1.520083423425557
iteration 800 / 2000: loss 1.5002372373288377
iteration 900 / 2000: loss 1.5352295535942588
iteration 1000 / 2000: loss 1.4202376005841388
iteration 1100 / 2000: loss 1.4045676360997463
iteration 1200 / 2000: loss 1.4420523917033683
iteration 1300 / 2000: loss 1.40279133845263
iteration 1400 / 2000: loss 1.357727650253214
iteration 1500 / 2000: loss 1.5724928126039976
iteration 1600 / 2000: loss 1.5295854369566435
iteration 1700 / 2000: loss 1.4108235797343034
iteration 1800 / 2000: loss 1.4966362537588063
iteration 1900 / 2000: loss 1.4778364297144924
54 Iteration
Learning rate: 0.001 Reg: 0.3 Batch size: 300 Val_acc: 0.496
```

```
iteration 0 / 2000: loss 2.302791300549903
iteration 100 / 2000: loss 1.8818459127211618
iteration 200 / 2000: loss 1.8098969430365355
iteration 300 / 2000: loss 1.7322375781152504
iteration 400 / 2000: loss 1.587905878491272
iteration 500 / 2000: loss 1.5691290929036075
iteration 600 / 2000: loss 1.5569490805603248
iteration 700 / 2000: loss 1.463420824478733
iteration 800 / 2000: loss 1.537724394631854
iteration 900 / 2000: loss 1.4850845972113362
iteration 1000 / 2000: loss 1.3975151111031774
iteration 1100 / 2000: loss 1.5311832360440347
iteration 1200 / 2000: loss 1.3468226163065922
iteration 1300 / 2000: loss 1.4180207739060033
iteration 1400 / 2000: loss 1.4693785403028585
iteration 1500 / 2000: loss 1.4203041666804175
iteration 1600 / 2000: loss 1.4097370793851531
iteration 1700 / 2000: loss 1.3932902554801623
iteration 1800 / 2000: loss 1.3904900573192442
iteration 1900 / 2000: loss 1.3579794078327738
55 Iteration
Learning rate: 0.001 Reg: 0.3 Batch size: 400 Val_acc: 0.516
```

```
iteration 0 / 2000: loss 2.3028265677745665
```

```

iteration 100 / 2000: loss 1.9067867699202747
iteration 200 / 2000: loss 1.80630805461237
iteration 300 / 2000: loss 1.6499624762551997
iteration 400 / 2000: loss 1.6054280699357935
iteration 500 / 2000: loss 1.5436173283422985
iteration 600 / 2000: loss 1.5597406449940332
iteration 700 / 2000: loss 1.5663366434989745
iteration 800 / 2000: loss 1.4356148160731899
iteration 900 / 2000: loss 1.5360461162045225
iteration 1000 / 2000: loss 1.4770121852983558
iteration 1100 / 2000: loss 1.569131386577806
iteration 1200 / 2000: loss 1.4629839681293995
iteration 1300 / 2000: loss 1.4394980687808296
iteration 1400 / 2000: loss 1.4571694392945151
iteration 1500 / 2000: loss 1.5311082492115566
iteration 1600 / 2000: loss 1.4668068034695139
iteration 1700 / 2000: loss 1.448869660808112
iteration 1800 / 2000: loss 1.368153515436651
iteration 1900 / 2000: loss 1.4640161694041585
60 Iteration
Learning rate: 0.001 Reg: 0.5 Batch size: 500 Val_acc: 0.485

```

```

Best Learning rate: 0.001
Best Reg: 0.3
Best Batch Size: 400
Best Val_acc: 0.516
Best num_iteration: 2000
Validation accuracy: 0.516

```

```

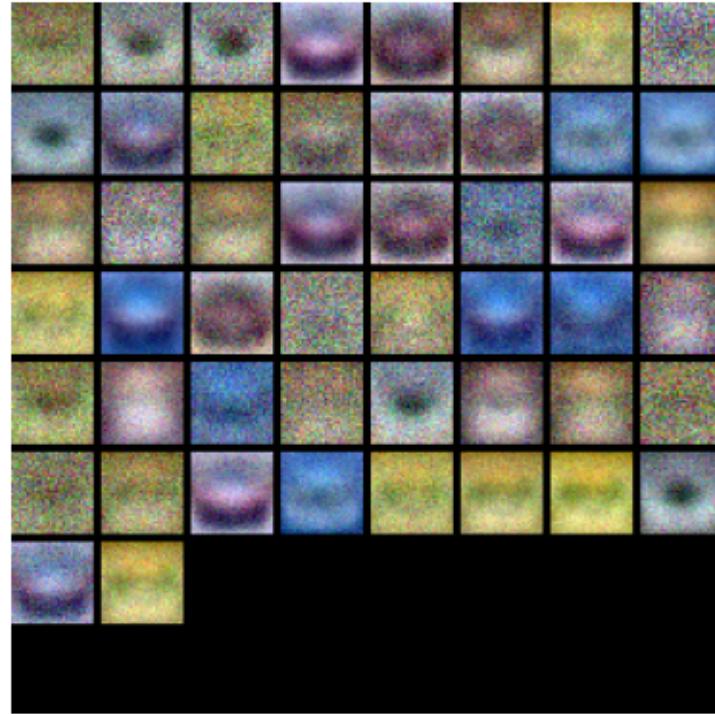
[ ]: from utils.vis_utils import visualize_grid

# Visualize the weights of the network

def show_net_weights(net):
    W1 = net.params['W1']
    W1 = W1.T.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
    plt.imshow(visualize_grid(W1, padding=3).astype('uint8'))
    plt.gca().axis('off')
    plt.show()

show_net_weights(subopt_net)
show_net_weights(best_net)

```



0.7 Question:

- (1) What differences do you see in the weights between the suboptimal net and the best net you arrived at?

0.8 Answer:

- (1) The weights of best net is much clearer than the suboptimal net and the boundary inside the weights image is also more obvious.
- (2) We can observe more details in best net. For example, in the suboptimal net, the image on 1st row and 4th col, let's call it (1,4) looks almost the same as (3,4). But in best net, we can find that actually (1,4) and (3,4) still look different. It means the weights of best net could help us obtain more imformation of the features.

0.9 Evaluate on test set

```
[ ]: test_acc = (best_net.predict(X_test) == y_test).mean()
      print('Test accuracy: ', test_acc)
```

Test accuracy: 0.518

neural_net

February 2, 2023

```
[ ]: # %load neural_net.py
import numpy as np
import matplotlib.pyplot as plt

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network. The net has an input dimension of
    D, a hidden layer dimension of H, and performs classification over C classes.
    We train the network with a softmax loss function and L2 regularization on the
    weight matrices. The network uses a ReLU nonlinearity after the first fully
    connected layer.

    In other words, the network has the following architecture:

    input - fully connected layer - ReLU - fully connected layer - softmax

    The outputs of the second fully-connected layer are the scores for each class.
    """

    def __init__(self, input_size, hidden_size, output_size, std=1e-4):
        """
        Initialize the model. Weights are initialized to small random values and
        biases are initialized to zero. Weights and biases are stored in the
        variable self.params, which is a dictionary with the following keys:
        W1: First layer weights; has shape (H, D)
        b1: First layer biases; has shape (H,)
        W2: Second layer weights; has shape (C, H)
        b2: Second layer biases; has shape (C,)

        Inputs:
        - input_size: The dimension D of the input data.
        - hidden_size: The number of neurons H in the hidden layer.
        - output_size: The number of classes C.
        """
        self.params = {}
```

```

    self.params['W1'] = std * np.random.randn(hidden_size, input_size)
    self.params['b1'] = np.zeros(hidden_size)
    self.params['W2'] = std * np.random.randn(output_size, hidden_size)
    self.params['b2'] = np.zeros(output_size)

def loss(self, X, y=None, reg=0.0):
    """
    Compute the loss and gradients for a two layer fully connected neural
    network.

    Inputs:
    - X: Input data of shape (N, D). Each X[i] is a training sample.
    - y: Vector of training labels. y[i] is the label for X[i], and each y[i] is
      an integer in the range 0 <= y[i] < C. This parameter is optional; if it
      is not passed then we only return scores, and if it is passed then we
      instead return the loss and gradients.
    - reg: Regularization strength.

    Returns:
    If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
    the score for class c on input X[i].
    If y is not None, instead return a tuple of:
    - loss: Loss (data loss and regularization loss) for this batch of training
      samples.
    - grads: Dictionary mapping parameter names to gradients of those parameters
      with respect to the loss function; has the same keys as self.params.
    """
    # Unpack variables from the params dictionary
    W1, b1 = self.params['W1'], self.params['b1']
    W2, b2 = self.params['W2'], self.params['b2']
    N, D = X.shape

    # Compute the forward pass
    scores = None

    # ===== #
    # YOUR CODE HERE:
    # Calculate the output scores of the neural network. The result
    # should be (N, C). As stated in the description for this class,
    # there should not be a ReLU layer after the second FC layer.
    # The output of the second FC layer is the output scores. Do not
    # use a for loop in your implementation.
    # ===== #
    X_after_first_layer = np.dot(X, W1.T) + b1 #broadcast
    X_after_ReLu = np.maximum(0, X_after_first_layer)

```

```

scores = np.dot(X_after_ReLu, W2.T) + b2
pass

# ===== #
# END YOUR CODE HERE
# ===== #

# If the targets are not given then jump out, we're done
if y is None:
    return scores

# Compute the loss
loss = None

# ===== #
# YOUR CODE HERE:
# Calculate the loss of the neural network. This includes the
# softmax loss and the L2 regularization for W1 and W2. Store the
# total loss in teh variable loss. Multiply the regularization
# loss by 0.5 (in addition to the factor reg).
# ===== #

N,C = scores.shape[0],scores.shape[1]
H,D = self.params['W1'].shape
scores = scores - np.max(scores, axis = 1, keepdims = True)

true_matrix = np.zeros((N,C))
true_matrix[np.arange(N), y] = 1
true_score_mat = np.multiply(true_matrix, scores)
true_score = np.sum(true_score_mat)

false_score = np.log( np.sum( np.exp(scores), axis = 1 ) ).sum()

loss = false_score - true_score
loss = loss / N + 0.5 * reg * np.sum(W1**2) + 0.5 * reg * np.sum(W2**2)
# scores is num_examples by num_classes
pass

# ===== #
# END YOUR CODE HERE
# ===== #

grads = {}

# ===== #
# YOUR CODE HERE:
# Implement the backward pass. Compute the derivatives of the
# weights and the biases. Store the results in the grads

```

```

#     dictionary. e.g., grads['W1'] should store the gradient for
#     W1, and be of the same size as W1.
# ===== #
d_scores = np.zeros((N,C))
d_scores = np.exp(scores) / np.exp(scores).sum(axis = 1, keepdims = True)
d_scores = (d_scores - true_matrix) / N
grads['W2'] = ( X_after_ReLu.T.dot(d_scores) + reg * W2.T ).T
grads['b2'] = np.sum(d_scores, axis = 0).T

d_input2 = np.zeros((H,D))
d_input2 = np.dot(d_scores, W2)
d_input2[X_after_ReLu < 1e-5] = 0
# print(d_input2.shape)
grads['b1'] = np.sum(d_input2, axis = 0).T
grads['W1'] = (X.T.dot(d_input2) + reg * W1.T).T
# print(grads['W1'].shape, grads['b1'].shape, grads['W2'].shape, □
↪grads['b2'].shape)

pass

# ===== #
# END YOUR CODE HERE
# ===== #

return loss, grads

```

```

def train(self, X, y, X_val, y_val,
          learning_rate=1e-3, learning_rate_decay=0.95,
          reg=1e-5, num_iters=100,
          batch_size=200, verbose=False):
    """

```

Train this neural network using stochastic gradient descent.

Inputs:

- *X: A numpy array of shape (N, D) giving training data.*
- *y: A numpy array f shape (N,) giving training labels; y[i] = c means that X[i] has label c, where 0 <= c < C.*
- *X_val: A numpy array of shape (N_val, D) giving validation data.*
- *y_val: A numpy array of shape (N_val,) giving validation labels.*
- *learning_rate: Scalar giving learning rate for optimization.*
- *learning_rate_decay: Scalar giving factor used to decay the learning rate after each epoch.*
- *reg: Scalar giving regularization strength.*
- *num_iters: Number of steps to take when optimizing.*
- *batch_size: Number of training examples to use per step.*
- *verbose: boolean; if true print progress during optimization.*

```

num_train = X.shape[0]
num_val = X_val.shape[0]
iterations_per_epoch = max(num_train / batch_size, 1)

# Use SGD to optimize the parameters in self.model
loss_history = []
train_acc_history = []
val_acc_history = []

for it in np.arange(num_iters):
    X_batch = None
    y_batch = None

    # ===== #
    # YOUR CODE HERE:
    # Create a minibatch by sampling batch_size samples randomly.
    # ===== #
    idx = np.random.choice(range(num_train), batch_size)
    X_batch = X[idx]
    y_batch = y[idx]
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    # Compute loss and gradients using the current minibatch
    loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
    loss_history.append(loss)

    # ===== #
    # YOUR CODE HERE:
    # Perform a gradient descent step using the minibatch to update
    # all parameters (i.e., W1, W2, b1, and b2).
    # ===== #
    self.params['W1'] = self.params['W1'] - learning_rate * grads['W1']
    self.params['W2'] = self.params['W2'] - learning_rate * grads['W2']
    self.params['b1'] = self.params['b1'] - learning_rate * grads['b1']
    self.params['b2'] = self.params['b2'] - learning_rate * grads['b2']
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    if verbose and it % 100 == 0:
        print('iteration {} / {}: loss {}'.format(it, num_iters, loss))

    # Every epoch, check train and val accuracy and decay learning rate.
    if it % iterations_per_epoch == 0:
        # Check accuracy

```

```

        train_acc = (self.predict(X_batch) == y_batch).mean()
        val_acc = (self.predict(X_val) == y_val).mean()
        train_acc_history.append(train_acc)
        val_acc_history.append(val_acc)

        # Decay learning rate
        learning_rate *= learning_rate_decay

    return {
        'loss_history': loss_history,
        'train_acc_history': train_acc_history,
        'val_acc_history': val_acc_history,
    }

def predict(self, X):
    """
    Use the trained weights of this two-layer network to predict labels for
    data points. For each data point we predict scores for each of the C
    classes, and assign each data point to the class with the highest score.

    Inputs:
    - X: A numpy array of shape (N, D) giving N D-dimensional data points to
      classify.

    Returns:
    - y_pred: A numpy array of shape (N,) giving predicted labels for each of
      the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
      to have class c, where 0 <= c < C.
    """
    y_pred = None
    # ===== #
    # YOUR CODE HERE:
    # Predict the class given the input data.
    # ===== #
    X_after_first_layer = np.dot(X, self.params['W1'].T) + self.params['b1'] ↴#broadcast
    X_after_ReLu = np.maximum(0, X_after_first_layer)
    scores = np.dot(X_after_ReLu, self.params['W2'].T) + self.params['b2']
    y_pred = np.argmax(scores, axis = 1)
    pass

    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return y_pred

```

fc_nets_ipynb

February 5, 2023

1 Fully connected networks

In the previous notebook, you implemented a simple two-layer neural network class. However, this class is not modular. If you wanted to change the number of layers, you would need to write a new loss and gradient function. If you wanted to optimize the network with different optimizers, you'd need to write new training functions. If you wanted to incorporate regularizations, you'd have to modify the loss and gradient function.

Instead of having to modify functions each time, for the rest of the class, we'll work in a more modular framework where we define forward and backward layers that calculate losses and gradients respectively. Since the forward and backward layers share intermediate values that are useful for calculating both the loss and the gradient, we'll also have these functions return "caches" which store useful intermediate values.

The goal is that through this modular design, we can build different sized neural networks for various applications.

In this HW #3, we'll define the basic architecture, and in HW #4, we'll build on this framework to implement different optimizers and regularizations (like BatchNorm and Dropout).

1.1 Modular layers

This notebook will build modular layers in the following manner. First, there will be a forward pass for a given layer with inputs (`x`) and return the output of that layer (`out`) as well as cached variables (`cache`) that will be used to calculate the gradient in the backward pass.

```
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:

```

def layer_backward(dout, cache):
    """
    Receive 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

```

[]: ## Import and setups

```

import time
import numpy as np
import matplotlib.pyplot as plt
from ndl.fc_net import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, □
    ↪eval_numerical_gradient_array
from utils.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.maximum(1e-8, np.abs(x) + np.abs(y))))

```

[]: # Load the (preprocessed) CIFAR10 data.

```

data = get_CIFAR10_data()
for k in data.keys():
    print('{}: {}'.format(k, data[k].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,)
```

1.2 Linear layers

In this section, we'll implement the forward and backward pass for the linear layers.

The linear layer forward pass is the function `affine_forward` in `nndl/layers.py` and the backward pass is `affine_backward`.

After you have implemented these, test your implementation by running the cell below.

1.2.1 Affine layer forward pass

Implement `affine_forward` and then test your code by running the following cell.

```
[ ]: # 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 1e-9.  
print('Testing affine_forward function:')print('difference: {}'.format(rel_error(out, correct_out)))
```

```
Testing affine_forward function:  
difference: 9.7698500479884e-10
```

1.2.2 Affine layer backward pass

Implement `affine_backward` and then test your code by running the following cell.

```
[ ]: # Test the affine_backward function  
  
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 1e-10
print('Testing affine_backward function:')
print('dx error: {}'.format(rel_error(dx_num, dx)))
print('dw error: {}'.format(rel_error(dw_num, dw)))
print('db error: {}'.format(rel_error(db_num, db)))

```

```

Testing affine_backward function:
dx error: 8.845861628974646e-10
dw error: 4.405537926313691e-10
db error: 1.531984065068008e-11

```

1.3 Activation layers

In this section you'll implement the ReLU activation.

1.3.1 ReLU forward pass

Implement the `relu_forward` function in `nndl/layers.py` and then test your code by running the following cell.

```

[ ]: # 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 around 1e-8
print('Testing relu_forward function:')
print('difference: {}'.format(rel_error(out, correct_out)))

```

```

Testing relu_forward function:
difference: 4.99999798022158e-08

```

1.3.2 ReLU backward pass

Implement the `relu_backward` function in `nndl/layers.py` and then test your code by running the following cell.

```
[ ]: 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 around 1e-12
print('Testing relu_backward function:')
print('dx error: {}'.format(rel_error(dx_num, dx)))
```

Testing `relu_backward` function:
dx error: 3.2756144884343895e-12

1.4 Combining the affine and ReLU layers

Often times, an affine layer will be followed by a ReLU layer. So let's make one that puts them together. Layers that are combined are stored in `nndl/layer_utils.py`.

1.4.1 Affine-ReLU layers

We've implemented `affine_relu_forward()` and `affine_relu_backward` in `nndl/layer_utils.py`. Take a look at them to make sure you understand what's going on. Then run the following cell to ensure its implemented correctly.

```
[ ]: from nndl.layer_utils import affine_relu_forward, affine_relu_backward

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)

print('Testing affine_relu_forward and affine_relu_backward:')
```

```

print('dx error: {}'.format(rel_error(dx_num, dx)))
print('dw error: {}'.format(rel_error(dw_num, dw)))
print('db error: {}'.format(rel_error(db_num, db)))

```

Testing affine_relu_forward and affine_relu_backward:

```

dx error: 1.0361713691913885e-10
dw error: 1.3235761825149517e-10
db error: 3.2756031511796114e-12

```

1.5 Softmax loss

You've already implemented it, so we have written it in `layers.py`. The following code will ensure they are working correctly.

```

[ ]: 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 2.3 and dx error should be 1e-8
print('\nTesting softmax_loss:')
print('loss: {}'.format(loss))
print('dx error: {}'.format(rel_error(dx_num, dx)))

```

Testing softmax_loss:

```

loss: 2.3024541667232965
dx error: 8.740178208266035e-09

```

1.6 Implementation of a two-layer NN

In `nndl/fc_net.py`, implement the class `TwoLayerNet` which uses the layers you made here. When you have finished, the following cell will test your implementation.

```

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

std = 1e-2
model = TwoLayerNet(input_dim=D, hidden_dims=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],
     [12.05769098, 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'

for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = {}'.format(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('{} relative error: {}'.format(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: 1.2236151215593397e-08
W2 relative error: 3.3429539606923665e-10
b1 relative error: 4.7288944058018464e-09
b2 relative error: 4.3291285233961314e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 3.1196421170137957e-07
W2 relative error: 7.976669724128995e-08
b1 relative error: 1.3467621663681023e-08
b2 relative error: 7.759087852670534e-10
b2 relative error: 7.759087852670534e-10

```

1.7 Solver

We will now use the utils Solver class to train these networks. Familiarize yourself with the API in `utils/solver.py`. After you have done so, declare an instance of a TwoLayerNet with 200 units and then train it with the Solver. Choose parameters so that your validation accuracy is at least 50%.

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

# ===== #
# YOUR CODE HERE:
# Declare an instance of a TwoLayerNet and then train
# it with the Solver. Choose hyperparameters so that your validation
# accuracy is at least 50%. We won't have you optimize this further
# since you did it in the previous notebook.
#
# ===== #
model = TwoLayerNet(hidden_dims = 200)
X_train, y_train = data['X_train'], data['y_train']
X_val, y_val = data['X_val'], data['y_val']
X_test, y_test = data['X_test'], data['y_test']
best_lr = 0.001
best_batch_size = 400
best_num_iter = 2000
best_reg = 0.3
mydata = {
    'X_train': X_train,# training data
    'y_train': y_train,# training labels
    'X_val': X_val,# validation data
    'y_val': y_val# validation labels
}
# model = TwoLayerNet(hidden_size=100, reg=best_reg)
```

```

solver = Solver(model, mydata,
                update_rule='sgd',
                optim_config={
                    'learning_rate': best_lr,
                },
                lr_decay=0.95,
                num_epochs=10, batch_size=best_batch_size,
                print_every=100)

solver.train()
pass

# ===== #
# END YOUR CODE HERE
# ===== #

```

```

(Iteration 1 / 1220) loss: 2.298761
(Epoch 0 / 10) train acc: 0.188000; val_acc: 0.144000
(Iteration 101 / 1220) loss: 1.786055
(Epoch 1 / 10) train acc: 0.381000; val_acc: 0.393000
(Iteration 201 / 1220) loss: 1.573889
(Epoch 2 / 10) train acc: 0.466000; val_acc: 0.436000
(Iteration 301 / 1220) loss: 1.546084
(Epoch 3 / 10) train acc: 0.441000; val_acc: 0.462000
(Iteration 401 / 1220) loss: 1.510700
(Epoch 4 / 10) train acc: 0.518000; val_acc: 0.465000
(Iteration 501 / 1220) loss: 1.478797
(Iteration 601 / 1220) loss: 1.429531
(Epoch 5 / 10) train acc: 0.524000; val_acc: 0.484000
(Iteration 701 / 1220) loss: 1.339706
(Epoch 6 / 10) train acc: 0.531000; val_acc: 0.498000
(Iteration 801 / 1220) loss: 1.340428
(Epoch 7 / 10) train acc: 0.503000; val_acc: 0.507000
(Iteration 901 / 1220) loss: 1.297518
(Epoch 8 / 10) train acc: 0.537000; val_acc: 0.484000
(Iteration 1001 / 1220) loss: 1.311881
(Epoch 9 / 10) train acc: 0.547000; val_acc: 0.517000
(Iteration 1101 / 1220) loss: 1.339683
(Iteration 1201 / 1220) loss: 1.318799
(Epoch 10 / 10) train acc: 0.543000; val_acc: 0.530000

```

[]: # Run this cell to visualize training loss and train / val accuracy

```

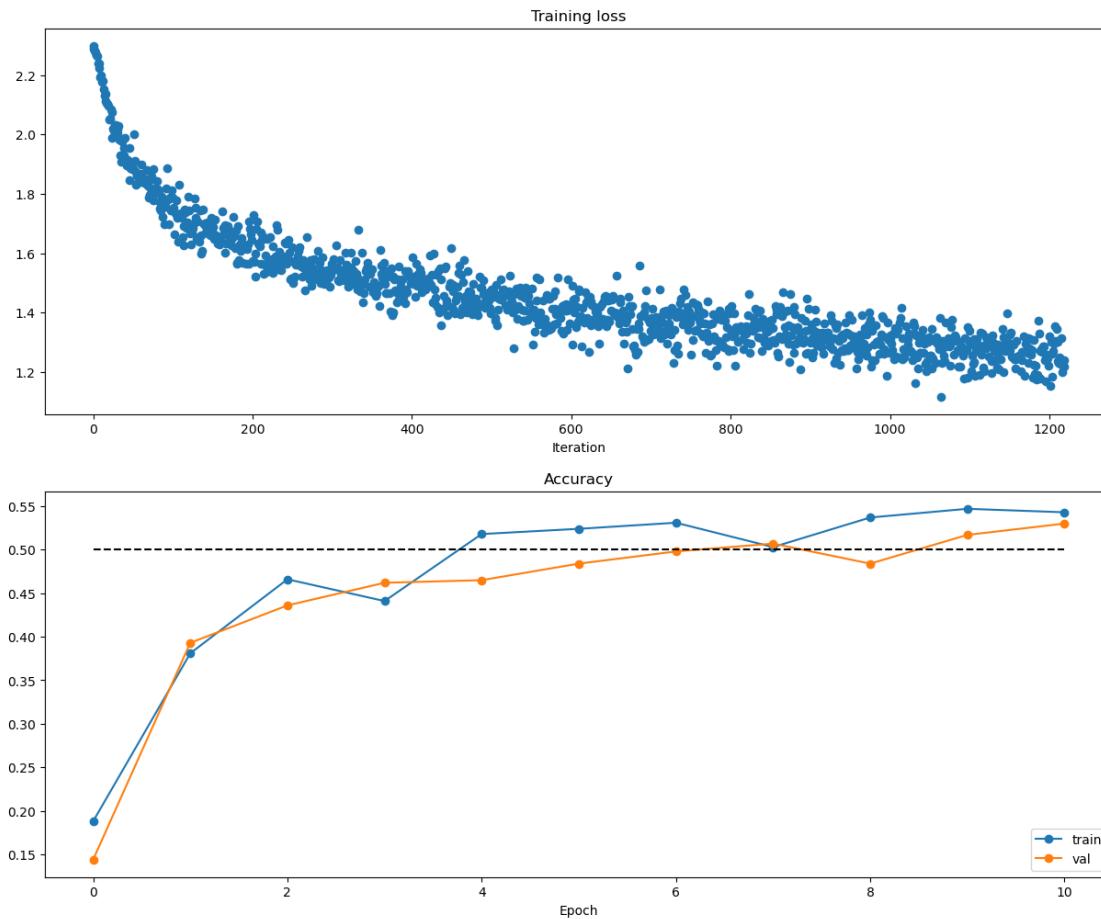
plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
plt.xlabel('Iteration')

```

```

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()

```



1.8 Multilayer Neural Network

Now, we implement a multi-layer neural network.

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

Implement the initialization, the forward pass, and the backward pass. There will be lines for batchnorm and dropout layers and caches; ignore these all for now. That'll be in HW #4.

```
[ ]: 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 = {}'.format(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: {}'.format(loss))

    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('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
```

```
Running check with reg = 0
Initial loss: 2.2386625318598483
W1 relative error: 4.591955823949023e-08
W2 relative error: 1.1142544919980604e-07
W3 relative error: 8.727606850444996e-08
b1 relative error: 1.5861303017086138e-09
b2 relative error: 3.470836247759264e-10
b3 relative error: 9.392121566378062e-11
Running check with reg = 3.14
Initial loss: 509.957466505932
W1 relative error: 4.957414338464139e-07
W2 relative error: 8.487616752528423e-06
W3 relative error: 3.0690488903966404e-07
b1 relative error: 7.902554834430548e-08
b2 relative error: 7.77362582871271e-08
b3 relative error: 1.3738787948178761e-08
```

```
[ ]: # Use the three layer neural network to overfit a small dataset.

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'],
}
```

```

##### !!!!!!
# Play around with the weight_scale and learning_rate so that you can overfit a
# small dataset.
# Your training accuracy should be 1.0 to receive full credit on this part.
weight_scale = 1e-2
learning_rate = 1e-4

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()

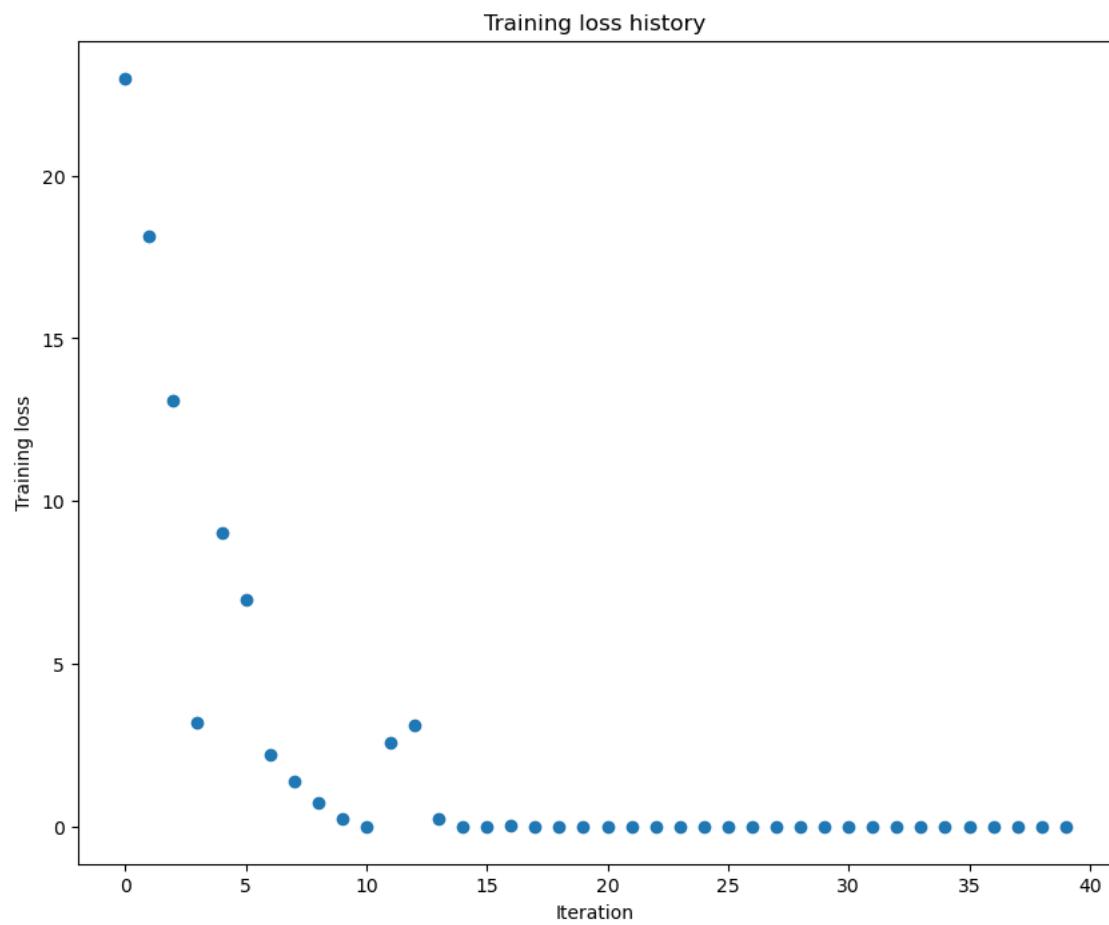
```

```

(Iteration 1 / 40) loss: 22.997194
(Epoch 0 / 20) train acc: 0.300000; val_acc: 0.130000
(Epoch 1 / 20) train acc: 0.500000; val_acc: 0.127000
(Epoch 2 / 20) train acc: 0.580000; val_acc: 0.130000
(Epoch 3 / 20) train acc: 0.760000; val_acc: 0.138000
(Epoch 4 / 20) train acc: 0.920000; val_acc: 0.163000
(Epoch 5 / 20) train acc: 0.940000; val_acc: 0.187000
(Iteration 11 / 40) loss: 0.000828
(Epoch 6 / 20) train acc: 0.880000; val_acc: 0.158000
(Epoch 7 / 20) train acc: 0.980000; val_acc: 0.179000
(Epoch 8 / 20) train acc: 0.980000; val_acc: 0.178000
(Epoch 9 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 10 / 20) train acc: 1.000000; val_acc: 0.177000
(Iteration 21 / 40) loss: 0.002572
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.176000
(Iteration 31 / 40) loss: 0.001593
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.178000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.177000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.178000

```

```
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.178000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.178000
```



```
[ ]:
```

layers

February 2, 2023

```
[ ]: # %load layers.py
import numpy as np
import pdb
```

```
def affine_forward(x, w, b):
```

```
    """
```

Computes the forward pass for an affine (fully-connected) layer.

*The input x has shape (N, d_1, \dots, d_k) and contains a minibatch of N examples, where each example $x[i]$ has shape (d_1, \dots, d_k) . We will reshape each input into a vector of dimension $D = d_1 * \dots * d_k$, and then transform it to an output vector of dimension M .*

Inputs:

- x : A numpy array containing input data, of shape (N, d_1, \dots, d_k)
- w : A numpy array of weights, of shape (D, M)
- b : A numpy array of biases, of shape $(M,)$

Returns a tuple of:

- out : output, of shape (N, M)
- $cache$: (x, w, b)

```
    """
```

```
# ===== #
# YOUR CODE HERE:
# Calculate the output of the forward pass. Notice the dimensions
# of w are  $D \times M$ , which is the transpose of what we did in earlier
# assignments.
# ===== #
num_inputs = x.shape[0]
input_shape = x.shape[1:]
input_size = np.prod(input_shape)
x_reshape = x.reshape(num_inputs, input_size)
out = np.dot(x_reshape, w) + b
```

```

pass

# ===== #
# END YOUR CODE HERE
# ===== #

cache = (x, w, b)
return out, cache


def affine_backward(dout, cache):
    """
    Computes the backward pass for an affine layer.

    Inputs:
    - dout: Upstream derivative, of shape (N, M)
    - cache: Tuple of:
        - x: Input data, of shape (N, d_1, ..., d_k)
        - w: Weights, of shape (D, M)

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """
    x, w, b = cache
    dx, dw, db = None, None, None

    # ===== #
    # YOUR CODE HERE:
    # Calculate the gradients for the backward pass.
    # ===== #

    # dout is N x M
    # dx should be N x d1 x ... x dk; it relates to dout through multiplication with w, which is D x M
    # dw should be D x M; it relates to dout through multiplication with x, which is N x D after reshaping
    # db should be M; it is just the sum over dout examples
    num_inputs = x.shape[0]
    input_shape = x.shape[1:]
    input_size = np.prod(input_shape)
    x_reshape = x.reshape(num_inputs, input_size)
    x_shape = x.shape
    db = np.sum(dout, axis=0)
    dx = np.dot(dout, w.T).reshape(x_shape)
    dw = np.dot(x_reshape.T, dout)

```

```

pass

# ===== #
# END YOUR CODE HERE
# ===== #

return dx, dw, db

def relu_forward(x):
    """
    Computes the forward pass for a layer of rectified linear units (ReLUs).

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # ===== #
    # YOUR CODE HERE:
    # Implement the ReLU forward pass.
    # ===== #
    out = np.maximum(x, 0)
    pass
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    cache = x
    return out, cache


def relu_backward(dout, cache):
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).

    Input:
    - dout: Upstream derivatives, of any shape
    - cache: Input x, of same shape as dout

    Returns:
    - dx: Gradient with respect to x
    """
    x = cache

    # ===== #

```

```

# YOUR CODE HERE:
# Implement the ReLU backward pass
# ===== #
x[x<0] = 0
x[x>0] = 1
dx = np.multiply(x,dout)
# ReLU directs linearly to those > 0
pass

# ===== #
# END YOUR CODE HERE
# ===== #

return dx

def softmax_loss(x, y):
    """
    Computes the loss and gradient for softmax classification.

    Inputs:
    - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
          for the ith input.
    - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
          0 <= y[i] < C

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    probs = np.exp(x - np.max(x, axis=1, keepdims=True))
    probs /= np.sum(probs, axis=1, keepdims=True)
    N = x.shape[0]
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx

```

fc_net

February 2, 2023

```
[ ]: # %load fc_net.py
import numpy as np

from .layers import *
from .layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network with ReLU nonlinearity and
    softmax loss that uses a modular layer design. We assume an input dimension
    of D, a hidden dimension of H, and perform classification over C classes.

    The architecture should be affine - relu - affine - softmax.

    Note that this class does not implement gradient descent; instead, it
    will interact with a separate Solver object that is responsible for running
    optimization.

    The learnable parameters of the model are stored in the dictionary
    self.params that maps parameter names to numpy arrays.

    """
    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
                 dropout=0, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dims: An integer giving the size of the hidden layer
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
            initialization of the weights.
        - reg: Scalar giving L2 regularization strength.
        """

```

```

self.params = {}
self.reg = reg

# ===== #
# YOUR CODE HERE:
# Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
# self.params['W2'], self.params['b1'] and self.params['b2']. The
# biases are initialized to zero and the weights are initialized
# so that each parameter has mean 0 and standard deviation weight_scale.
# The dimensions of W1 should be (input_dim, hidden_dim) and the
# dimensions of W2 should be (hidden_dims, num_classes)
# ===== #
self.params['W1'] = np.random.randn(input_dim, hidden_dims) * weight_scale
self.params['W2'] = np.random.randn(hidden_dims, num_classes) * weight_scale
self.params['b1'] = np.zeros((1,hidden_dims))
self.params['b2'] = np.zeros((1,num_classes))
pass

# ===== #
# END YOUR CODE HERE
# ===== #

def loss(self, X, y=None):
    """
    Compute loss and gradient for a minibatch of data.

    Inputs:
    - X: Array of input data of shape (N, d_1, ..., d_k)
    - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
    """

    if y is None:
        # If y is None, then run a test-time forward pass of the model and return:
        # scores: Array of shape (N, C) giving classification scores, where
        # scores[i, c] is the classification score for X[i] and class c.
    else:
        # If y is not None, then run a training-time forward and backward pass and
        # return a tuple of:
        # - loss: Scalar value giving the loss
        # - grads: Dictionary with the same keys as self.params, mapping parameter
        # names to gradients of the loss with respect to those parameters.
    """
    scores = None

    # ===== #
    # YOUR CODE HERE:
    # Implement the forward pass of the two-layer neural network. Store
    # the class scores as the variable 'scores'. Be sure to use the layers

```

```

#     you prior implemented.
# ===== #
W1, W2 = self.params['W1'], self.params['W2']
b1, b2 = self.params['b1'], self.params['b2']
output_first_layer, cache_first = affine_relu_forward(X, W1, b1)
scores, cache_sec = affine_forward(output_first_layer, W2,b2)
pass
# ===== #
# END YOUR CODE HERE
# ===== #

# If y is None then we are in test mode so just return scores
if y is None:
    return scores

loss, grads = 0, {}
# ===== #
# YOUR CODE HERE:
# Implement the backward pass of the two-layer neural net. Store
# the loss as the variable 'loss' and store the gradients in the
# 'grads' dictionary. For the grads dictionary, grads['W1'] holds
# the gradient for W1, grads['b1'] holds the gradient for b1, etc.
# i.e., grads[k] holds the gradient for self.params[k].
#
# Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
# for each W. Be sure to include the 0.5 multiplying factor to
# match our implementation.
#
# And be sure to use the layers you prior implemented.
# ===== #
loss, dout = softmax_loss(scores ,y)
d_first_hidden, dW2, db2 = affine_backward(dout, cache_sec)
dx, dW1, db1 = affine_relu_backward(d_first_hidden, cache_first)
dW1 = dW1 + self.reg * W1
dW2 = dW2 + self.reg * W2
grads['W1'], grads['W2'] = dW1, dW2
grads['b1'], grads['b2'] = db1, db2
loss = loss + 0.5 * np.sum( self.reg * W1 * W1 ) + 0.5 * np.sum( self.reg * ↴W2 * W2 )

pass
# ===== #
# END YOUR CODE HERE
# ===== #

return loss, grads

```

```

class FullyConnectedNet(object):
    """
    A fully-connected neural network with an arbitrary number of hidden layers,
    ReLU nonlinearities, and a softmax loss function. This will also implement
    dropout and batch normalization as options. For a network with L layers,
    the architecture will be

    {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax

    where batch normalization and dropout are optional, and the {...} block is
    repeated L - 1 times.

    Similar to the TwoLayerNet above, learnable parameters are stored in the
    self.params dictionary and will be learned using the Solver class.
    """

    def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
                 dropout=0, use_batchnorm=False, reg=0.0,
                 weight_scale=1e-2, dtype=np.float32, seed=None):
        """
        Initialize a new FullyConnectedNet.

        Inputs:
        - hidden_dims: A list of integers giving the size of each hidden layer.
        - input_dim: An integer giving the size of the input.
        - num_classes: An integer giving the number of classes to classify.
        - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
          the network should not use dropout at all.
        - use_batchnorm: Whether or not the network should use batch normalization.
        - reg: Scalar giving L2 regularization strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - dtype: A numpy datatype object; all computations will be performed using
          this datatype. float32 is faster but less accurate, so you should use
          float64 for numeric gradient checking.
        - seed: If not None, then pass this random seed to the dropout layers. This
          will make the dropout layers deterministic so we can gradient check the
          model.
        """
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout > 0
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}

```

```

# ===== #
# YOUR CODE HERE:
# Initialize all parameters of the network in the self.params dictionary.
# The weights and biases of layer 1 are W1 and b1; and in general the
# weights and biases of layer i are Wi and bi. The
# biases are initialized to zero and the weights are initialized
# so that each parameter has mean 0 and standard deviation weight_scale.
# ===== #
self.params['W1'] = np.random.randn(input_dim, hidden_dims[0]) * weight_scale
self.params['b1'] = np.zeros(hidden_dims[0])

for i in range(self.num_layers - 2):
    j = i + 2
    weights = str('W') + str(j)
    bias = str('b') + str(j)
    # print(weights,bias)
    self.params[weights] = np.random.randn(hidden_dims[j-2], hidden_dims[j-1]) * weight_scale
    self.params[bias] = np.zeros((1, hidden_dims[j-1]))

weight_last = str('W') + str(self.num_layers)
bias_last = str('b') + str(self.num_layers)
self.params[weight_last] = np.random.randn(hidden_dims[-1], num_classes)
self.params[bias_last] = np.zeros(num_classes)
pass

# ===== #
# END YOUR CODE HERE
# ===== #

# When using dropout we need to pass a dropout_param dictionary to each
# dropout layer so that the layer knows the dropout probability and the mode
# (train / test). You can pass the same dropout_param to each dropout layer.
self.dropout_param = {}
if self.use_dropout:
    self.dropout_param = {'mode': 'train', 'p': dropout}
    if seed is not None:
        self.dropout_param['seed'] = seed

# With batch normalization we need to keep track of running means and
# variances, so we need to pass a special bn_param object to each batch
# normalization layer. You should pass self.bn_params[0] to the forward pass
# of the first batch normalization layer, self.bn_params[1] to the forward

```

```

# pass of the second batch normalization layer, etc.
self.bn_params = []
if self.use_batchnorm:
    self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - ↴1)]
    ↴1)

# Cast all parameters to the correct datatype
for k, v in self.params.items():
    self.params[k] = v.astype(dtype)

def loss(self, X, y=None):
    """
    Compute loss and gradient for the fully-connected net.

    Input / output: Same as TwoLayerNet above.
    """
    X = X.astype(self.dtype)
    mode = 'test' if y is None else 'train'

    # Set train/test mode for batchnorm params and dropout param since they
    # behave differently during training and testing.
    if self.dropout_param is not None:
        self.dropout_param['mode'] = mode
    if self.use_batchnorm:
        for bn_param in self.bn_params:
            bn_param[mode] = mode

    scores = None
    # ===== #
    # YOUR CODE HERE:
    # Implement the forward pass of the FC net and store the output
    # scores as the variable "scores".
    # ===== #
    output = {}
    output[0] = X
    cache = {}
    for i in range(self.num_layers):
        j = i + 1
        weights = str('W') + str(j)
        bias = str('b') + str(j)
        W_j = self.params[weights]
        b_j = self.params[bias]

        if j == self.num_layers:
            output[j], cache[i] = affine_forward(output[i], W_j, b_j)
        else:

```

```

        output[j], cache[i] = affine_relu_forward(output[i], W_j, b_j)

scores = output[self.num_layers]
pass

# ===== #
# END YOUR CODE HERE
# ===== #

# If test mode return early
if mode == 'test':
    return scores

loss, grads = 0.0, {}
# ===== #
# YOUR CODE HERE:
#   Implement the backwards pass of the FC net and store the gradients
#   in the grads dict, so that grads[k] is the gradient of self.params[k]
#   Be sure your L2 regularization includes a 0.5 factor.
# ===== #
loss, dx = softmax_loss(scores, y)
reg_loss = 0
for i in range(self.num_layers):
    j = i + 1
    weights = str('W') + str(j)
    W_j = self.params[weights]
    reg_loss = reg_loss + 0.5 * self.reg * np.sum(W_j * W_j)
loss = loss + reg_loss

# compute grads
doutput = {}
hidden_dims = self.num_layers - 1
weight_last = str('W') + str(self.num_layers)
bias_last = str('b') + str(self.num_layers)
W_last, b_last = self.params[weight_last], self.params[bias_last]
doutput[hidden_dims], grads[weight_last], grads[bias_last] = ↵
affine_backward(dx, cache[hidden_dims])
grads[weight_last] = grads[weight_last] + self.reg * self.
params[weight_last]

for i in range(hidden_dims):
    weight = str('W') + str(hidden_dims - i)
    bias = str('b') + str(hidden_dims - i)
    doutput[hidden_dims - i - 1], grads[weight], grads[bias] = ↵
affine_relu_backward(doutput[hidden_dims - i], cache[hidden_dims - i - 1])
grads[weight] = grads[weight] + self.reg * self.params[weight]

```

```
pass

# ===== #
# END YOUR CODE HERE
# ===== #
return loss, grads
```