

# ECE C147/247 HW4 Q1: Optimization for Fully Connected Networks

In this notebook, we will implement different optimization rules for gradient descent. We have provided starter code; however, you will need to copy and paste your code from your implementation of the modular fully connected nets in HW #3 to build upon this.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, and their layer structure. This also includes `nndl.fc_net`, `nndl.layers`, and `nndl.layer_utils`.

```
In [2]: ## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.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))))
```

```
In [3]: # 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,)
```

## Building upon your HW #3 implementation

Copy and paste the following functions from your HW #3 implementation of a modular FC net:

- `affine_forward` in `nndl/layers.py`
- `affine_backward` in `nndl/layers.py`
- `relu_forward` in `nndl/layers.py`
- `relu_backward` in `nndl/layers.py`
- `affine_relu_forward` in `nndl/layer_utils.py`
- `affine_relu_backward` in `nndl/layer_utils.py`
- The `FullyConnectedNet` class in `nndl/fc_net.py`

## Test all functions you copy and pasted

```
In [4]: from nndl.layer_tests import *

affine_forward_test(); print('\n')
affine_backward_test(); print('\n')
relu_forward_test(); print('\n')
relu_backward_test(); print('\n')
affine_relu_test(); print('\n')
fc_net_test()

If affine_forward function is working, difference should be less than 1e-9:
difference: 9.769849468192957e-10

If affine_backward is working, error should be less than 1e-9:
dx error: 4.6163154161927085e-10
dw error: 6.707494642048358e-11
db error: 6.942320637660904e-12

If relu_forward function is working, difference should be around 1e-8:
difference: 4.99999798022158e-08

If relu_backward function is working, error should be less than 1e-9:
dx error: 3.2756279540991835e-12

If affine_forward, relu_forward and affine_relu_backward are working, error should be less than 1e-9:
dx error: 5.524297922807348e-11
dw error: 3.631734171907549e-10
db error: 3.2756041823344864e-12

Running check with reg = 0
Initial loss: 2.303392659269333
W1 relative error: 2.865557200343655e-05
W2 relative error: 4.989612955448939e-07
W3 relative error: 8.249390118912997e-08
b1 relative error: 3.68580394215281e-07
b2 relative error: 3.4440948223621555e-09
b3 relative error: 1.4614235824389521e-10
Running check with reg = 3e-4
Initial loss: 5.744233115944262
W1 relative error: 4.160769835235782e-09
W2 relative error: 2.7898679841333852e-08
W3 relative error: 1.0
b1 relative error: 2.234842945592196e-06
b2 relative error: 1.1539627622763397e-08
b3 relative error: 2.6543776080429063e-10
```

## Training a larger model

In general, proceeding with vanilla stochastic gradient descent to optimize models may be fraught with problems and limitations, as discussed in class. Thus, we implement optimizers that improve on SGD.

## SGD + momentum

In the following section, implement SGD with momentum. Read the `nndl/optim.py` API, and be sure you understand it. After, implement `sgd_momentum` in `nndl/optim.py`. Test your implementation of `sgd_momentum` by running the cell below.

```
In [5]: from nndl.optim import sgd_momentum

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

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

expected_next_w = np.asarray([
    [0.406, 0.207389, 0.27417895, 0.34096842, 0.40775789],
    [0.47454737, 0.5413684, 0.60812632, 0.67491579, 0.74170526],
    [0.80849474, 0.87528421, 0.94207368, 1.00886316, 1.07565263],
    [1.14244211, 1.20923158, 1.27602105, 1.34281053, 1.4096 ]])
expected_velocity = np.asarray([
    [0.5406, 0.55475789, 0.56891579, 0.58307368, 0.59723158],
    [0.61138947, 0.62554737, 0.63970526, 0.65386316, 0.66802105],
    [0.68217895, 0.69633684, 0.71049474, 0.72465263, 0.73881053],
    [0.75296842, 0.76712632, 0.78128421, 0.79544211, 0.8096 ]])

print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))

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

## SGD + Nesterov momentum

Implement `sgd_nesterov_momentum` in `nndl/optim.py`.

```
In [6]: from nndl.optim import sgd_nesterov_momentum

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

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

expected_next_w = np.asarray([
    [0.08714, 0.15246105, 0.21778211, 0.28310316, 0.34842421],
    [0.41374526, 0.47906632, 0.54438737, 0.60970842, 0.67502947],
    [0.74035053, 0.80567158, 0.87099263, 0.93631368, 1.00163474],
    [1.06695579, 1.13227684, 1.19759789, 1.26291895, 1.3282447 ]])
expected_velocity = np.asarray([
    [0.5406, 0.55475789, 0.56891579, 0.58307368, 0.59723158],
    [0.61138947, 0.62554737, 0.63970526, 0.65386316, 0.66802105],
    [0.68217895, 0.69633684, 0.71049474, 0.72465263, 0.73881053],
    [0.75296842, 0.76712632, 0.78128421, 0.79544211, 0.8096 ]])

print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))

next_w error: 1.0875186845081027e-08
velocity error: 4.269287743278663e-09
```

## Evaluating SGD, SGD+Momentum, and SGD+NesterovMomentum

Run the following cell to train a 6-layer FC net with SGD, SGD+momentum, and SGD+Nesterov momentum. You should see that SGD+momentum achieves a better loss than SGD, and that SGD+Nesterov momentum achieves a slightly better loss (and training accuracy) than SGD+momentum.

```
In [7]: num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = []

for update_rule in ['sgd', 'sgd_momentum', 'sgd_nesterov_momentum']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100], weight_scale=5e-2)

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

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

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

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

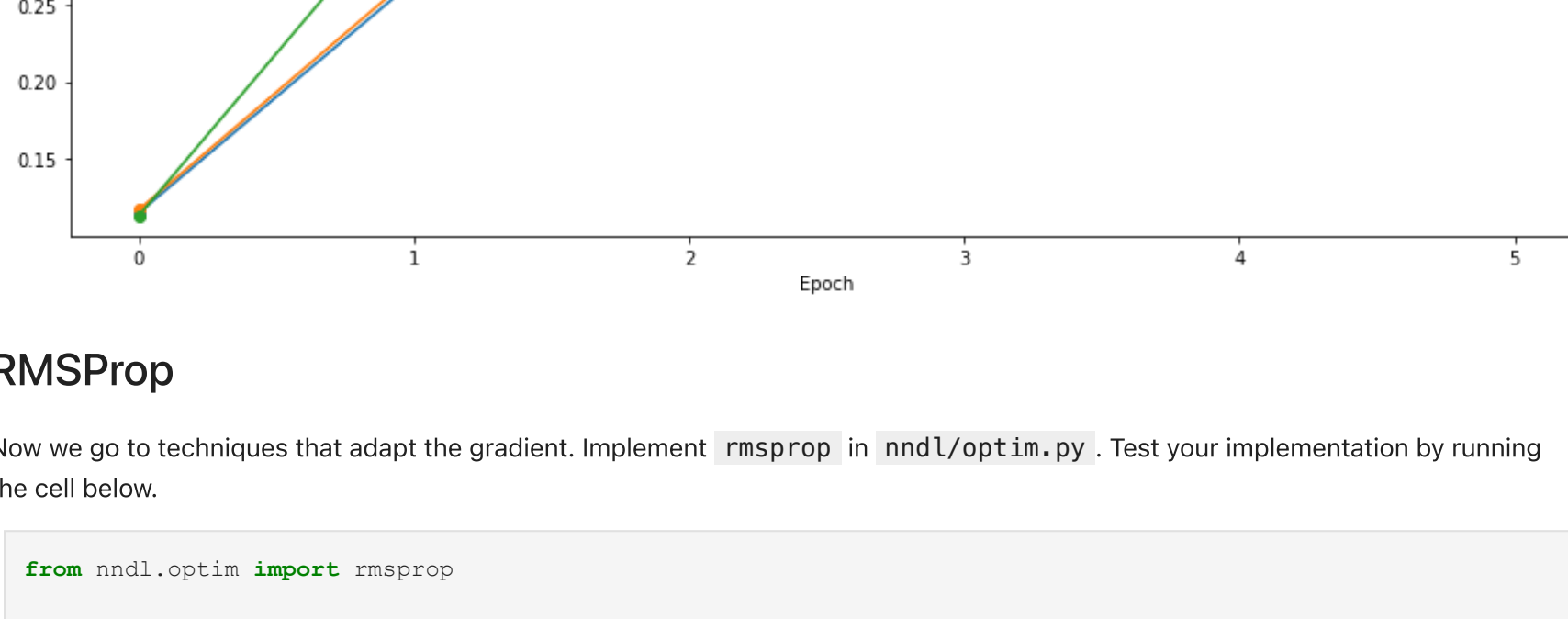
for update_rule, solver in solvers.items():
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

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

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

for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

Optimizing with sgd  
Optimizing with sgd\_momentum  
Optimizing with sgd\_nesterov\_momentum



Training loss

Training accuracy

Validation accuracy

## RMSProp

Now we go to techniques that adapt the gradient. Implement `rmsprop` in `nndl/optim.py`. Test your implementation by running the cell below.

```
In [8]: from nndl.optim import rmsprop

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

config = {'learning_rate': 1e-2, 'a': a}
next_w, _ = rmsprop(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
    [-0.132737, -0.08078555, -0.02891884, 0.02316247, 0.07515741],
    [0.12716641, 0.17918792, 0.23122175, 0.28326742, 0.33524471],
    [0.38739248, 0.43947102, 0.49155973, 0.54365823, 0.59576619]])
expected_cache = np.asarray([
    [0.5976, 0.6126277, 0.62771108, 0.64284931, 0.65804321],
    [0.67329252, 0.68859723, 0.70389734, 0.71927285, 0.73484377],
    [0.75037008, 0.7659518, 0.78158892, 0.79728144, 0.81302936],
    [0.82883269, 0.84469141, 0.86060554, 0.87657507, 0.8926 ]])

print('next_w error: {}'.format(rel_error(next_w, next_w)))
print('cache error: {}'.format(rel_error(expected_cache, config['a'])))

next_w error: 9.502645228984295e-08
cache error: 2.6477955807156126e-09
```

## Adaptive moments

Now, implement `adam` in `nndl/optim.py`. Test your implementation by running the cell below.

```
In [9]: # Test Adam implementation; you should see errors around 1e-7 or less

from nndl.optim import adam

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
a = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'v': v, 'a': a, 't': 5}
next_w, _ = adam(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
    [-0.1380274, -0.08544591, -0.03286534, 0.01971428, 0.07229291],
    [0.1248705, 0.17744702, 0.23002243, 0.28293667, 0.33516969],
    [0.38774145, 0.44031118, 0.49288082, 0.54544852, 0.59801459]])
expected_a = np.asarray([
    [0.69966, 0.6890382, 0.67851319, 0.66794809, 0.65738853],
    [0.64683452, 0.63628604, 0.6257431, 0.61520571, 0.60467385],
    [0.59414753, 0.58362676, 0.57311582, 0.56260183, 0.55209767],
    [0.54153906, 0.53110598, 0.52061845, 0.51013645, 0.49966 ]])
expected_v = np.asarray([
    [0.48, 0.49947368, 0.51894737, 0.53842105, 0.55789474],
    [0.57736842, 0.59684211, 0.61631579, 0.63578947, 0.65526316],
    [0.67473684, 0.69421053, 0.71368421, 0.73315789, 0.75263158],
    [0.77210526, 0.79157895, 0.81105263, 0.83052632, 0.85 ]])

print('next_w error: {}'.format(rel_error(next_w, next_w)))
print('a error: {}'.format(rel_error(expected_a, config['a'])))
print('v error: {}'.format(rel_error(expected_v, config['v'])))

next_w error: 1.1395691798535431e-07
a error: 4.208314038113071e-09
v error: 4.414963193114416e-09
```

## Comparing SGD, SGD+NesterovMomentum, RMSProp, and Adam

The following code will compare optimization with SGD, Momentum, Nesterov Momentum, RMSProp and Adam. In our code, we find that RMSProp, Adam, and SGD + Nesterov Momentum achieve approximately the same training error after a few training epochs.

```
In [10]: learning_rates = {'rmsprop': 2e-4, 'adam': 1e-3}

for update_rule in ['adam', 'rmsprop']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, data,
                    num_epochs=10, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': learning_rates[update_rule]
                    },
                    verbose=False)
    solvers[update_rule] = solver
    solver.train()
    print

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

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

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

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

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

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

for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

Optimizing with adam  
Optimizing with rmsprop



Training loss

Training accuracy

Validation accuracy

## Easier optimization

In the following cell, we'll train a 4+ layer neural network having 500 units in each hidden layer with the different optimizers, and find that it is far easier to get up to 50+% performance on CIFAR-10. After we implement batchnorm and dropout, we'll ask you to get 55+% on CIFAR-10.

```
In [11]: optimizer = 'adam'
best_model = None

layer_dims = [500, 500, 500]
weight_scale = 0.01
lr_decay = 0.9

model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
                          use_batchnorm=True)

solver = Solver(model, data,
                num_epochs=10, batch_size=100,
                update_rule=optimizer,
                optim_config={
                    'learning_rate': learning_rate,
                    'lr_decay': lr_decay,
                    'verbose': True, print_every=50
                })

solver.train()
```

(Iteration 1 / 4900) loss: 2.324108  
(Epoch 0 / 10) train acc: 0.259000; val\_acc: 0.236000  
(Iteration 51 / 4900) loss: 1.710812  
(Iteration 101 / 4900) loss: 1.569633  
(Iteration 151 / 4900) loss: 1.704779  
(Iteration 201 / 4900) loss: 1.553350  
(Iteration 251 / 4900) loss: 1.625926  
(Iteration 301 / 4900) loss: 1.565025  
(Iteration 351 / 4900) loss: 1.365191  
(Iteration 401 / 4900) loss: 1.479973  
(Iteration 451 / 4900) loss: 1.471955  
(Epoch 1 / 10) train acc: 0.469000; val\_acc: 0.480000  
(Iteration 501 / 4900) loss: 1.354020  
(Iteration 551 / 4900) loss: 1.429286  
(Iteration 601 / 4900) loss: 1.349315  
(Iteration 651 / 4900) loss: 1.021812  
(Iteration 701 / 4900) loss: 1.365892  
(Iteration 751 / 4900) loss: 1.391804  
(Iteration 801 / 4900) loss: 1.291474  
(Iteration 851 / 4900) loss: 1.355819  
(Iteration 901 / 4900) loss: 1.389257  
(Iteration 951 / 4900) loss: 1.358157  
(Epoch 2 / 10) train acc: 0.587000; val\_acc: 0.509000  
(Iteration 1001 / 4900) loss: 1.190631  
(Iteration 1051 / 4900) loss: 1.337119  
(Iteration 1101 / 4900) loss: 1.015182  
(Iteration 1151 / 4900) loss: 1.194532  
(Iteration 1201 / 4900) loss: 1.087486  
(Iteration 1251 / 4900) loss: 1.280915  
(Iteration 1301 / 4900) loss: 1.267049  
(Iteration 1351 / 4900) loss: 1.082300  
(Iteration 1401 / 4900) loss: 0.998963  
(Iteration 1451 / 4900) loss: 0.952096  
(Epoch 3 / 10) train acc: 0.603000; val\_acc: 0.516000  
(Iteration 1501 / 4900) loss: 1.180900  
(Iteration 1551 / 4900) loss: 1.246343  
(Iteration 1601 / 4900) loss: 1.180023  
(Iteration 1651 / 4900) loss: 1.222134  
(Iteration 1701 / 4900) loss: 0.939548  
(Iteration 1751 / 4900) loss: 1.044596  
(Iteration 1801 / 4900) loss: 0.926830  
(Iteration 1851 / 4900) loss: 1.047942  
(Iteration 1901 / 4900) loss: 0.980795  
(Iteration 1951 / 4900) loss: 1.044814  
(Epoch 4 / 10) train acc: 0.616000; val\_acc: 0.533000  
(Iteration 2001 / 4900) loss: 1.002781  
(Iteration 2051 / 4900) loss: 0.906200  
(Iteration 2101 / 4900) loss: 1.015182  
(Iteration 2151 / 4900) loss: 0.796182  
(Iteration 2201 / 4900) loss: 0.831100  
(Iteration 2251 / 4900) loss: 0.943659  
(Iteration 2301 / 4900) loss: 0.874367  
(Iteration 2351 / 4900) loss: 1.139650  
(Iteration 2401 / 4900) loss: 0.897629  
(Epoch 5 / 10) train acc: 0.675000; val\_acc: 0.543000  
(Iteration 2451 / 4900) loss: 1.015501  
(Iteration 2501 / 4900) loss: 0.861091  
(Iteration 2551 / 4900) loss: 0.900259  
(Iteration 2601 / 4900) loss: 1.006299  
(Iteration 2651 / 4900) loss: 0.866825  
(Iteration 2701 / 4900) loss: 0.864173  
(Iteration 2751 / 4900) loss: 0.882573  
(Iteration 2801 / 4900) loss: 0.896280  
(Iteration 2851 / 4900) loss: 0.950030  
(Iteration 2901 / 4900) loss: 0.898759  
(Epoch 6 / 10) train acc: 0.714000; val\_acc: 0.547000  
(Iteration 2951 / 4900) loss: 0.630894  
(Iteration 3001 / 4900) loss: 0.776758  
(Iteration 3051 / 4900) loss: 0.717995  
(Iteration 3101 / 4900) loss: 0.637762  
(Iteration 3151 / 4900) loss: 0.797481  
(Iteration 3201 / 4900) loss: 0.789366  
(Iteration 3251 / 4900) loss: 0.792586  
(Iteration 3301 / 4900) loss: 0.918864  
(Iteration 3351 / 4900) loss: 0.820589  
(Iteration 3401 / 4900) loss: 0.740379  
(Epoch 7 / 10) train acc: 0.739000; val\_acc: 0.581000  
(Iteration 3451 / 4900) loss: 0.877669  
(Iteration 3501 / 4900) loss: 0.937993  
(Iteration 3551 / 4900) loss: 0.782571  
(Iteration 3601 / 4900) loss: 0.855775  
(Iteration 3651 / 4900) loss: 0.653582  
(Iteration 3701 / 4900) loss: 0.745576  
(Iteration 3751 / 4900) loss: 0.787957  
(Iteration 3801 / 4900) loss: 0.702431  
(Iteration 3851 / 4900) loss: 0.762737  
(Iteration 3901 / 4900) loss: 0.668990  
(Epoch 8 / 10) train acc: 0.740000; val\_acc: 0.539000  
(Iteration 3951 / 4900) loss: 0.585486  
(Iteration 4001 / 4900) loss: 0.665097  
(Iteration 4051 / 4900) loss: 0.682409  
(Iteration 4101 / 4900) loss: 0.759852  
(Iteration 4151 / 4900) loss: 0.621568  
(Iteration 4201 / 4900) loss: 0.445623  
(Iteration 4251 / 4900) loss: 0.703892  
(Iteration 4301 / 4900) loss: 0.804312  
(Iteration 4351 / 4900) loss: 0.650961  
(Iteration 4401 / 4900) loss: 0.528363  
(Epoch 9 / 10) train acc: 0.801000; val\_acc: 0.573000  
(Iteration 4451 / 4900) loss: 0.626053  
(Iteration 4501 / 4900) loss: 0.512020  
(Iteration 4551 / 4900) loss: 0.460116  
(Iteration 4601 / 4900) loss: 0.631010  
(Iteration 4651 / 4900) loss: 0.762133  
(Iteration 4701 / 4900) loss: 0.537167  
(Iteration 4751 / 4900) loss: 0.638203  
(Iteration 4801 / 4900) loss: 0.516493  
(Iteration 4851 / 4900) loss: 0.517716  
(Epoch 10 / 10) train acc: 0.830000; val\_acc: 0.558000

```
In [12]: y_test_pred = np.argmax(model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(model.loss(data['X_val']), axis=1)
print('Validation set accuracy: {}'.format(np.mean(y_val_pred == data['y_val'])))
print('Test set accuracy: {}'.format(np.mean(y_test_pred == data['y_test'])))
```

Validation set accuracy: 0.584  
Test set accuracy: 0.553

In [ ]:





# ECE C147/247 HW4 Q2: Batch Normalization

In this notebook, you will implement the batch normalization layers of a neural network to increase its performance. Please review the details of batch normalization from the lecture notes.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, and their layer structure. This also includes `nnDL.fc_net`, `nnDL.layers`, and `nnDL.layer_utils`.

```
In [1]: ## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nnDL.fc_net import *
from nnDL.layers 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))))

In [2]: # 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,)
```

## Batchnorm forward pass

Implement the training time batchnorm forward pass, `batchnorm_forward`, in `nnDL/Layers.py`. After that, test your implementation by running the following cell.

```
In [3]: # Check the training-time forward pass by checking means and variances
# of features both before and after batch normalization

# Simulate the forward pass for a two-layer network
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print(' means: ', a.mean(axis=0))
print(' stds: ', a.std(axis=0))

# Means should be close to zero and stds close to one
print('After batch normalization (gamma=1, beta=0)')
a_norm, _ = batchnorm_forward(a, np.ones(D3), np.zeros(D3), {'mode': 'train'})
print(' mean: ', a_norm.mean(axis=0))
print(' std: ', a_norm.std(axis=0))

# Now means should be close to beta and stds close to gamma
gamma = np.asarray([1.0, 2.0, 3.0])
beta = np.asarray([11.0, 12.0, 13.0])
a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print('After batch normalization (nontrivial gamma, beta)')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

Before batch normalization:
means: [49.25890874 15.68042077 10.77425698]
stds: [32.74515657 29.35640593 37.42221779]
After batch normalization (gamma=1, beta=0)
mean: [-2.54796184e-16 -4.71844785e-17 -3.72618603e-17]
std: [1. 0.99999999 1. ]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds: [1. 1.99999999 2.99999999]
```

Implement the testing time batchnorm forward pass, `batchnorm_forward`, in `nnDL/Layers.py`. After that, test your implementation by running the following cell.

```
In [4]: # Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
for t in np.arange(50):
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    batchnorm_forward(a, gamma, beta, bn_param)
    bn_param['mode'] = 'test'
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After batch normalization (test-time):')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

After batch normalization (test-time):
means: [ 0.11609137 -0.09502236  0.07165287]
stds: [ 1.05963918  0.92295591  0.99428131]
```

## Batchnorm backward pass

Implement the backward pass for the batchnorm layer, `batchnorm_backward` in `nnDL/Layers.py`. Check your implementation by running the following cell.

```
In [5]: # Gradient check batchnorm backward pass

N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

bn_param = {'mode': 'train'}
fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
fg = lambda a: batchnorm_forward(x, gamma, beta, bn_param)[0]
fb = lambda b: batchnorm_forward(x, gamma, beta, bn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma, dout)
db_num = eval_numerical_gradient_array(fb, beta, dout)

_, dgamma, dbeta = batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dgamma))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  5.586527822853507e-10
dgamma error:  2.6632508082038723e-11
dbeta error:  4.977243498503734e-11
```

## Implement a fully connected neural network with batchnorm layers

Modify the `FullyConnectedNet()` class in `nnDL/fc_net.py` to incorporate batchnorm layers. You will need to modify the class in the following areas:

- (1) The gammas and betas need to be initialized to 1's and 0's respectively in `__init__`.
- (2) The `batchnorm_forward` layer needs to be inserted between each affine and relu layer (except in the output layer) in a forward pass computation in `loss`. You may find it helpful to write an `affine_batchnorm_relu()` layer in `nnDL/layer_utils.py` although this is not necessary.
- (3) The `batchnorm_backward` layer has to be appropriately inserted when calculating gradients.

After you have done the appropriate modifications, check your implementation by running the following cell.

Note, while the relative error for W3 should be small, as we backprop gradients more, you may find the relative error increases. Our relative error for W1 is on the order of 1e-4.

```
In [6]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

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

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

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        dx_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
    if reg == 0: print('\n')
```

```
Running check with reg = 0
Initial loss: 2.3214776413153135
W1 relative error: 1.701456185897653e-05
W2 relative error: 5.150159402515397e-06
W3 relative error: 5.277463191023614e-10
b1 relative error: 0.0022204516003654358
b2 relative error: 5.551115123125783e-09
b3 relative error: 1.17055603130984187e-10
beta1 relative error: 4.104526065485512e-09
beta2 relative error: 1.0109375534584169e-08
gamma1 relative error: 3.2168407209000194e-09
gamma2 relative error: 5.788857349186422e-09
```

```
Running check with reg = 3.14
Initial loss: 5.961655036876346
W1 relative error: 3.338186289789542e-05
W2 relative error: 1.2190849287679875e-06
W3 relative error: 1.0
b1 relative error: 5.551115123125783e-09
b2 relative error: 4.440892098500626e-08
b3 relative error: 2.847632039326928e-10
beta1 relative error: 7.358215383341911e-09
beta2 relative error: 5.391805549094429e-09
gamma1 relative error: 7.616191542818951e-09
gamma2 relative error: 1.1586508381957641e-08
```

## Training a deep fully connected network with batch normalization.

To see if batchnorm helps, let's train a deep neural network with and without batch normalization.

```
In [7]: # Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

bn_solver = Solver(bn_model, small_data,
                  num_epochs=10, batch_size=50,
                  update_rule='adam',
                  optim_config={
                      'learning_rate': 1e-3,
                  },
                  verbose=True, print_every=200)
bn_solver.train()

solver = Solver(model, small_data,
                num_epochs=10, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=200)
solver.train()
```

```
(Iteration 1 / 200) loss: 2.294970
(Epoch 0 / 10) train acc: 0.157000; val_acc: 0.134000
(Epoch 1 / 10) train acc: 0.357000; val_acc: 0.272000
(Epoch 2 / 10) train acc: 0.407000; val_acc: 0.309000
(Epoch 3 / 10) train acc: 0.477000; val_acc: 0.338000
(Epoch 4 / 10) train acc: 0.539000; val_acc: 0.323000
(Epoch 5 / 10) train acc: 0.628000; val_acc: 0.322000
(Epoch 6 / 10) train acc: 0.669000; val_acc: 0.309000
(Epoch 7 / 10) train acc: 0.680000; val_acc: 0.309000
(Epoch 8 / 10) train acc: 0.668000; val_acc: 0.308000
(Epoch 9 / 10) train acc: 0.726000; val_acc: 0.334000
(Epoch 10 / 10) train acc: 0.795000; val_acc: 0.327000
(Iteration 1 / 200) loss: 2.302627
(Epoch 0 / 10) train acc: 0.163000; val_acc: 0.149000
(Epoch 1 / 10) train acc: 0.216000; val_acc: 0.201000
(Epoch 2 / 10) train acc: 0.304000; val_acc: 0.249000
(Epoch 3 / 10) train acc: 0.320000; val_acc: 0.279000
(Epoch 4 / 10) train acc: 0.350000; val_acc: 0.246000
(Epoch 5 / 10) train acc: 0.462000; val_acc: 0.334000
(Epoch 6 / 10) train acc: 0.469000; val_acc: 0.298000
(Epoch 7 / 10) train acc: 0.546000; val_acc: 0.320000
(Epoch 8 / 10) train acc: 0.569000; val_acc: 0.322000
(Epoch 9 / 10) train acc: 0.605000; val_acc: 0.336000
(Epoch 10 / 10) train acc: 0.660000; val_acc: 0.354000
```

```
In [9]: # Try training a very deep net with batchnorm
hidden_dims = [50, 50, 50, 50, 50, 50]

num_train = 1000
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'],
}

bn_solvers = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

    bn_solver = Solver(bn_model, small_data,
                      num_epochs=10, batch_size=50,
                      update_rule='adam',
                      optim_config={
                          'learning_rate': 1e-3,
                      },
                      verbose=False, print_every=200)
    bn_solver.train()

    solver = Solver(model, small_data,
                    num_epochs=10, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 1e-3,
                    },
                    verbose=False, print_every=200)
    solvers[weight_scale] = solver
```

```
Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20
/Users/jack_tseng/Documents/UCLA/Courses/winter_2021/c247/hw/hw4/assignment/hw4-code/nnDL/layers.py:426: RuntimeWarning: divide by zero encountered in log
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20
```

```
In [11]: # Plot results of best weight scale experiment
best_train_accs, bn_best_train_accs = [], []
best_val_accs, bn_best_val_accs = [], []
final_train_loss, bn_final_train_loss = [], []

for ws in weight_scales:
    best_train_accs.append(max(solvers[ws].train_acc_history))
    bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))

    best_val_accs.append(max(solvers[ws].val_acc_history))
    bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))

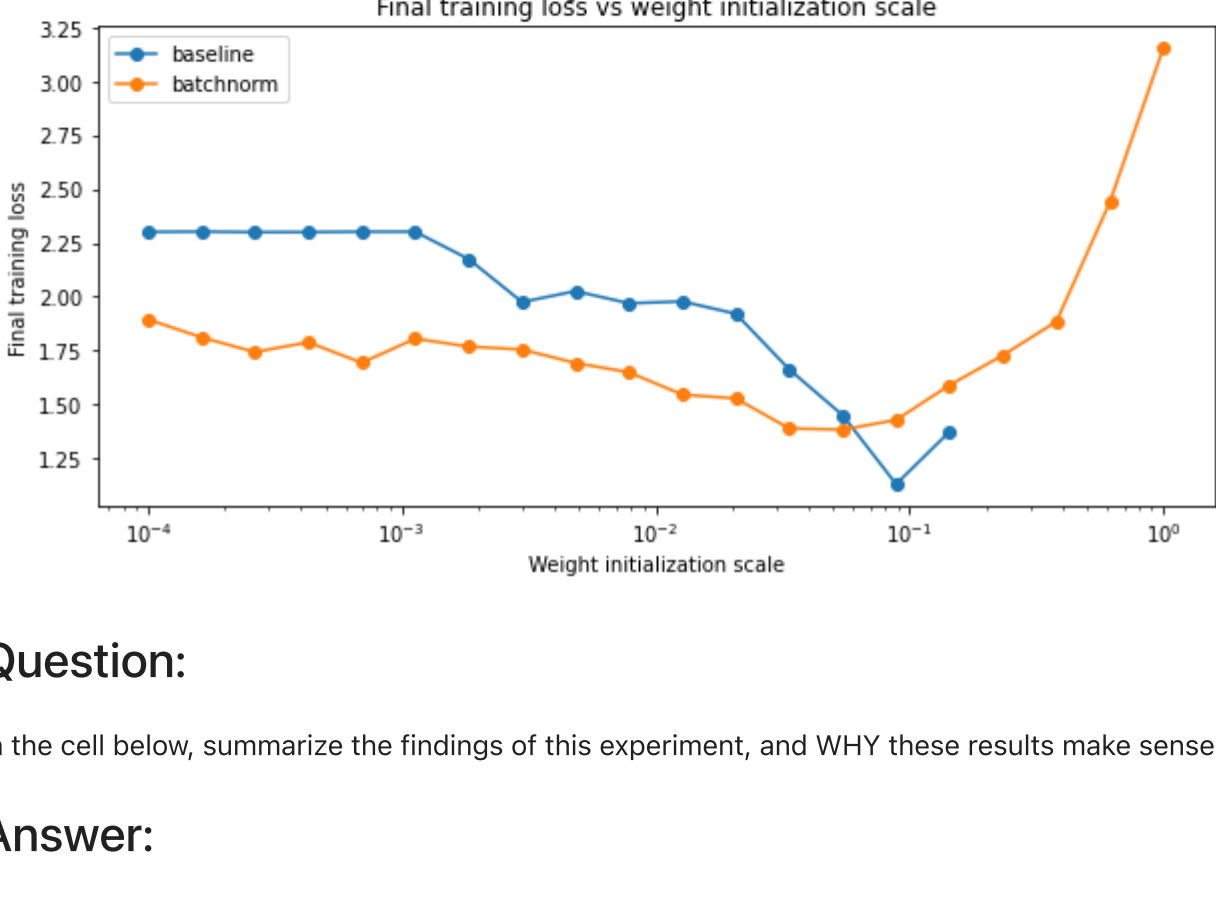
    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
plt.legend()

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()

plt.gcf().set_size_inches(10, 15)
plt.show()
```



## Question:

In the cell below, summarize the findings of this experiment, and WHY these results make sense.

## Answer:

Batchnorm would help stabilize the input data, also keep each backpropagations less effect by previous one. This would decrease the number of the required epoch. As we can see in the graph, while the epoch is small, batchnorm technique has better performance (~30% better). However, the benefits from batchnorm might vanish while the epoch is large.













## 1. Optimizer

```
def sgd_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

    # ===== #
    # YOUR CODE HERE:
    # Implement the momentum update formula. Return the updated weights
    # as next_w, and the updated velocity as v.
    # ===== #
    v = (config['momentum'] * v) - (config['learning_rate'] * dw)
    next_w = w + v
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    config['velocity'] = v

    return next_w, config
```

a.

```
def sgd_nesterov_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with Nesterov momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

    # ===== #
    # YOUR CODE HERE:
    # Implement the momentum update formula. Return the updated weights
    # as next_w, and the updated velocity as v.
    # ===== #
    v_old = v
    v = config['momentum'] * v_old - config['learning_rate'] * dw
    next_w = w + v + config['momentum'] * (v - v_old)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    config['velocity'] = v
```

b.

```
    return next_w, config
```

```

def rmsprop(w, dw, config=None):
    """
    Uses the RMSProp update rule, which uses a moving average of squared gradient
    values to set adaptive per-parameter learning rates.

    config format:
    - learning_rate: Scalar learning rate.
    - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
      gradient cache.
    - epsilon: Small scalar used for smoothing to avoid dividing by zero.
    - beta: Moving average of second moments of gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('decay_rate', 0.99)
    config.setdefault('epsilon', 1e-8)
    config.setdefault('a', np.zeros_like(w))

    next_w = None

    # ===== #
    # YOUR CODE HERE:
    # Implement RMSProp. Store the next value of w as next_w. You need
    # to also store in config['a'] the moving average of the second
    # moment gradients, so they can be used for future gradients. Concretely,
    # config['a'] corresponds to "a" in the lecture notes.
    # ===== #
    config['a'] = config['decay_rate'] * config['a'] + (1-config['decay_rate']) * np.multiply(dw, dw)
    next_w = w - np.multiply(config['learning_rate'] / (np.sqrt(config['a'] + config['epsilon'])), dw)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return next_w, config

```

C.

## 2. Batch Normalization

### a. Batch Forward

```

minibatch_mean = np.mean(x, axis=0)
minibatch_var = np.var(x, axis=0)
x_normalize = (x - minibatch_mean) / np.sqrt(minibatch_var + eps)
out = gamma * x_normalize + beta

running_mean = momentum * running_mean + (1 - momentum) * minibatch_mean
running_var = momentum * running_var + (1 - momentum) * minibatch_var
bn_param['running_mean'] = running_mean
bn_param['running_var'] = running_var

cache = {
    'minibatch_var': minibatch_var,
    'x_centralize': (x - minibatch_mean),
    'x_normalize': x_normalize,
    'gamma': gamma,
    'eps': eps
}
# ===== #
# END YOUR CODE HERE
# ===== #

elif mode == 'test':
    # ===== #
    # YOUR CODE HERE:
    # Calculate the testing time normalized activation. Normalize using
    # the running mean and variance, and then scale and shift appropriately.
    # Store the output as 'out'.
    # ===== #

    out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta

```

i.

### b. Batch Backward:



```

"""
dx, dgamma, dbeta = None, None, None

# ===== #
# YOUR CODE HERE:
# Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
# ===== #
N = dout.shape[0]
minibatch_var = cache.get('minibatch_var')
x_centralize = cache.get('x_centralize')
x_normalize = cache.get('x_normalize')
gamma = cache.get('gamma')
eps = cache.get('eps')

# calculate dx
dxhat = dout * gamma
dxmu1 = dxhat / np.sqrt(minibatch_var + eps)
sqrt_var = np.sqrt(minibatch_var + eps)
dsqrt_var = -np.sum(dxhat * x_centralize, axis=0) / (sqrt_var**2)
dvar = dsqrt_var * 0.5 / sqrt_var
dxmu2 = 2 * x_centralize * dvar * np.ones_like(dout) / N
dx1 = dxmu1 + dxmu2
dx2 = -np.sum(dx1, axis=0) * np.ones_like(dout) / N
dx = dx1 + dx2

# calculate dbeta and dgamma
dbeta = np.sum(dout, axis=0)
dgamma = np.sum(dout * x_normalize, axis=0)
# ===== #
# END YOUR CODE HERE
# ===== #

return dx, dgamma, dbeta

```

i.

### 3. Dropout

#### a. Dropout forward:

```

out = None

if mode == 'train':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout forward pass during training time.
    # Store the masked and scaled activations in out, and store the
    # dropout mask as the variable mask.
    # ===== #

    mask = (np.random.rand(x.shape[0], x.shape[1]) < p) / p
    out = np.multiply(x, mask)

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

elif mode == 'test':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout forward pass during test time.
    # ===== #

    mask = (np.random.random_sample(x.shape) >= p) / (1 - p)
    out = x * mask

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

    cache = (dropout_param, mask)

```

i.

#### b. Dropout Backward:

```

"""
dropout_param, mask = cache
mode = dropout_param['mode']

dx = None
if mode == 'train':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout backward pass during training time.
    # ===== #

    dx = dout * mask

    # ===== #
    # END YOUR CODE HERE
    # ===== #
elif mode == 'test':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout backward pass during test time.
    # ===== #

    dx = dout

    # ===== #
    # END YOUR CODE HERE
    # ===== #
return dx

```

i.