

Optimization

February 10, 2021

0.1 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.

CS231n 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`. As in prior assignments, we thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu).

```
[30]: ## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, \u
    ↪ eval_numerical_gradient_array
from cs231n.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))))
```

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

```
%reload_ext autoreload
```

```
[31]: # 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,)
```

0.2 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`

0.2.1 Test all functions you copy and pasted

```
[32]: 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: $1.3892139318961866e-10$
dw error: $2.5826575989869555e-09$
db error: $5.037883801368915e-12$

If `relu_forward` function is working, difference should be around $1e-8$:

difference: 4.999999798022158e-08

If relu_forward function is working, error should be less than 1e-9:
dx error: 3.2756086976798748e-12

If affine_relu_forward and affine_relu_backward are working, error should be less than 1e-9::
dx error: 2.753999848546374e-10
dw error: 4.591118402758519e-10
db error: 1.8928931789554515e-11

Running check with reg = 0
Initial loss: 2.301460950673409
W1 relative error: 1.5750832914729826e-06
W2 relative error: 8.413606662497265e-07
W3 relative error: 6.656190343813889e-07
b1 relative error: 4.2922263639427574e-08
b2 relative error: 1.8795965523873468e-09
b3 relative error: 1.2769957885744794e-10
Running check with reg = 3.14
Initial loss: 5.935472085131298
W1 relative error: 6.461256444442613e-08
W2 relative error: 2.0244961122371758e-08
W3 relative error: 1.0
b1 relative error: 5.900922224200075e-09
b2 relative error: 6.578576872297275e-09
b3 relative error: 2.8747933318374744e-10

1 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.

1.1 SGD + momentum

In the following section, implement SGD with momentum. Read the `nndl/optim.py` API, which is provided by CS231n, 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.

```
[33]: 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([
    [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
    [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
    [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
    [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096    ]])
expected_velocity = np.asarray([
    [ 0.5406,      0.55475789,  0.56891579,  0.58307368,  0.59723158],
    [ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
    [ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
    [ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096    ]])

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

```

1.2 SGD + Nesterov momentum

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

```

[34]: from ndl.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,  1.32824    ]])
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

```

1.3 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.

```

[35]: 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, 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)

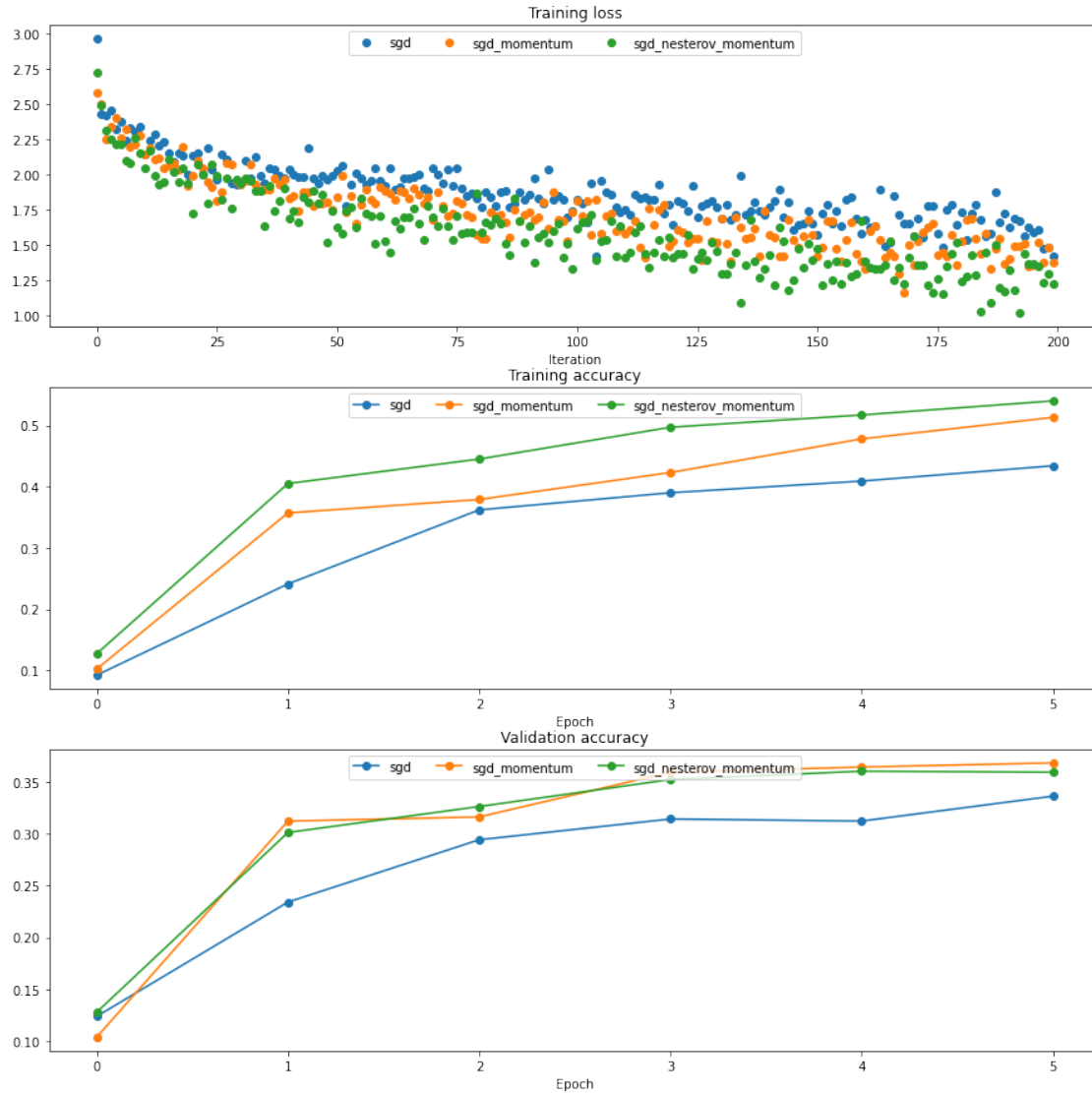
    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    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`



1.4 RMSProp

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

```
[36]: 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, -0.02881884, 0.02316247, 0.07515774],
    [ 0.12716641, 0.17918792, 0.23122175, 0.28326742, 0.33532447],
    [ 0.38739248, 0.43947102, 0.49155973, 0.54365823, 0.59576619]])
expected_cache = np.asarray([
    [ 0.5976, 0.6126277, 0.6277108, 0.64284931, 0.65804321],
    [ 0.67329252, 0.68859723, 0.70395734, 0.71937285, 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(expected_next_w, next_w)))
print('cache error: {}'.format(rel_error(expected_cache, config['a'])))

```

```

next_w error: 9.502645229894295e-08
cache error: 2.6477955807156126e-09

```

1.5 Adaptive moments

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

```

[37]: # 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, 0.0722929],
    [ 0.1248705, 0.17744702, 0.23002243, 0.28259667, 0.33516969],
    [ 0.38774145, 0.44031188, 0.49288093, 0.54544852, 0.59801459]])
expected_a = np.asarray([
    [ 0.69966, 0.68908382, 0.67851319, 0.66794809, 0.65738853],
    [ 0.64683452, 0.63628604, 0.6257431, 0.61520571, 0.60467385],
    [ 0.59414753, 0.58362676, 0.57311152, 0.56260183, 0.55209767],
    [ 0.54159906, 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(expected_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: 4.214963193114416e-09

```

1.6 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.

```

[38]: 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, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, 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)

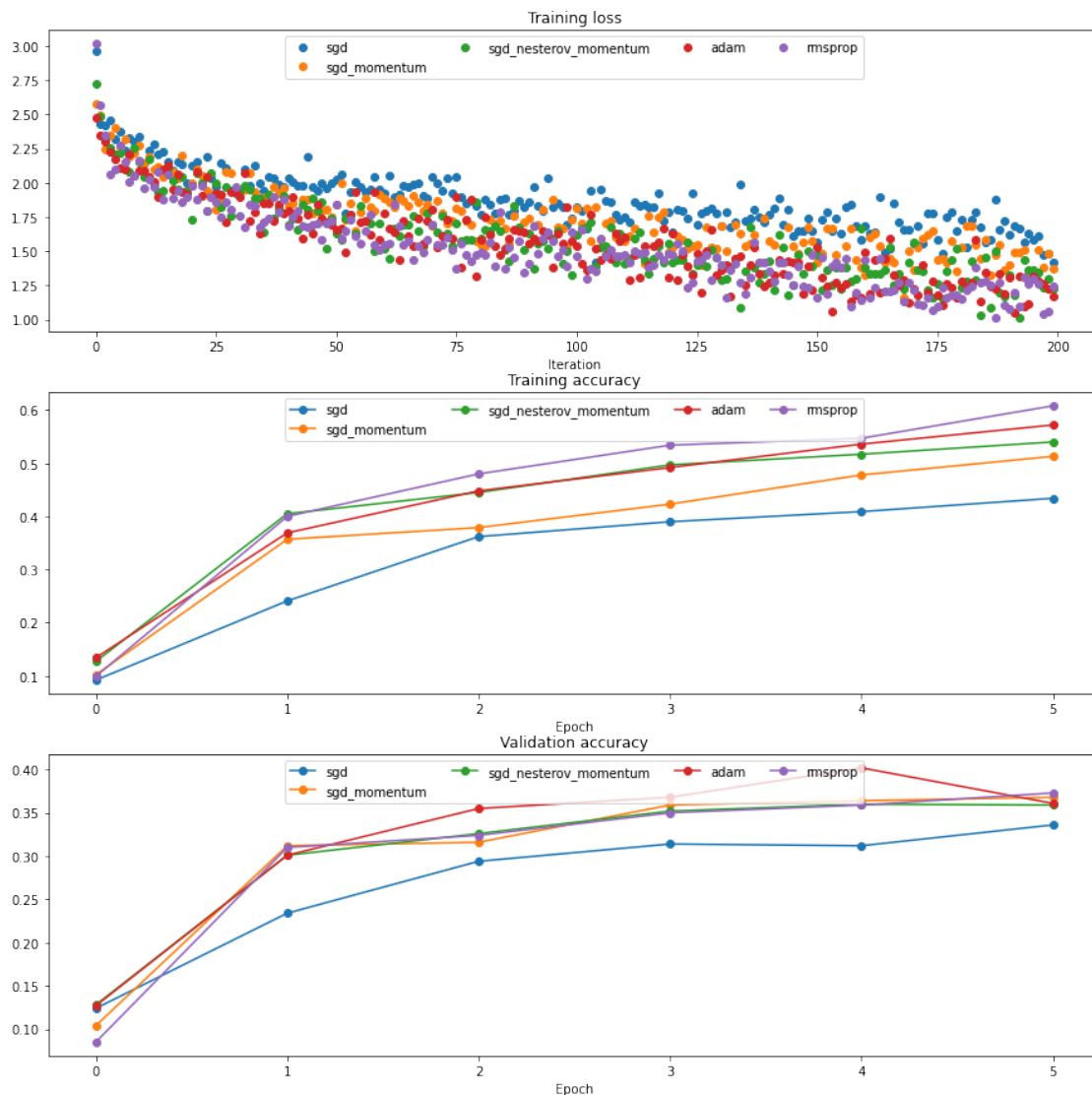
plt.subplot(3, 1, 3)
plt.plot(solver.val_acc_history, '-o', label=update_rule)

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

```

Optimizing with adam

Optimizing with rmsprop



1.7 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.

```
[39]: optimizer = 'adam'
      best_model = None

      layer_dims = [500, 500, 500]
      weight_scale = 0.01
      learning_rate = 1e-3
      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: 2.303304
(Epoch 0 / 10) train acc: 0.196000; val_acc: 0.195000
(Iteration 51 / 4900) loss: 1.753733
(Iteration 101 / 4900) loss: 1.607837
(Iteration 151 / 4900) loss: 1.594427
(Iteration 201 / 4900) loss: 1.406384
(Iteration 251 / 4900) loss: 1.577677
(Iteration 301 / 4900) loss: 1.243411
(Iteration 351 / 4900) loss: 1.488788
(Iteration 401 / 4900) loss: 1.485148
(Iteration 451 / 4900) loss: 1.407836
(Epoch 1 / 10) train acc: 0.495000; val_acc: 0.482000
(Iteration 501 / 4900) loss: 1.326885
(Iteration 551 / 4900) loss: 1.337456
(Iteration 601 / 4900) loss: 1.202338
(Iteration 651 / 4900) loss: 1.364537
(Iteration 701 / 4900) loss: 1.212084
(Iteration 751 / 4900) loss: 1.611742
```

(Iteration 801 / 4900) loss: 1.471870
(Iteration 851 / 4900) loss: 1.325310
(Iteration 901 / 4900) loss: 1.194848
(Iteration 951 / 4900) loss: 1.607586
(Epoch 2 / 10) train acc: 0.543000; val_acc: 0.498000
(Iteration 1001 / 4900) loss: 1.281553
(Iteration 1051 / 4900) loss: 1.118752
(Iteration 1101 / 4900) loss: 1.073722
(Iteration 1151 / 4900) loss: 1.161873
(Iteration 1201 / 4900) loss: 0.962861
(Iteration 1251 / 4900) loss: 1.199029
(Iteration 1301 / 4900) loss: 1.183931
(Iteration 1351 / 4900) loss: 1.216622
(Iteration 1401 / 4900) loss: 1.131508
(Iteration 1451 / 4900) loss: 1.102764
(Epoch 3 / 10) train acc: 0.562000; val_acc: 0.518000
(Iteration 1501 / 4900) loss: 1.299038
(Iteration 1551 / 4900) loss: 1.068179
(Iteration 1601 / 4900) loss: 1.080086
(Iteration 1651 / 4900) loss: 1.127397
(Iteration 1701 / 4900) loss: 1.158260
(Iteration 1751 / 4900) loss: 0.860723
(Iteration 1801 / 4900) loss: 1.066489
(Iteration 1851 / 4900) loss: 1.051771
(Iteration 1901 / 4900) loss: 1.355291
(Iteration 1951 / 4900) loss: 1.294773
(Epoch 4 / 10) train acc: 0.619000; val_acc: 0.535000
(Iteration 2001 / 4900) loss: 1.121633
(Iteration 2051 / 4900) loss: 0.874532
(Iteration 2101 / 4900) loss: 1.060572
(Iteration 2151 / 4900) loss: 1.132936
(Iteration 2201 / 4900) loss: 0.885685
(Iteration 2251 / 4900) loss: 1.052590
(Iteration 2301 / 4900) loss: 0.976331
(Iteration 2351 / 4900) loss: 0.881012
(Iteration 2401 / 4900) loss: 1.137029
(Epoch 5 / 10) train acc: 0.656000; val_acc: 0.535000
(Iteration 2451 / 4900) loss: 0.936625
(Iteration 2501 / 4900) loss: 1.047290
(Iteration 2551 / 4900) loss: 0.869810
(Iteration 2601 / 4900) loss: 1.006359
(Iteration 2651 / 4900) loss: 0.925071
(Iteration 2701 / 4900) loss: 0.926250
(Iteration 2751 / 4900) loss: 0.889879
(Iteration 2801 / 4900) loss: 0.882390
(Iteration 2851 / 4900) loss: 0.968382
(Iteration 2901 / 4900) loss: 0.806344
(Epoch 6 / 10) train acc: 0.688000; val_acc: 0.565000

```

(Iteration 2951 / 4900) loss: 0.699684
(Iteration 3001 / 4900) loss: 0.859054
(Iteration 3051 / 4900) loss: 0.686731
(Iteration 3101 / 4900) loss: 0.998542
(Iteration 3151 / 4900) loss: 0.878873
(Iteration 3201 / 4900) loss: 0.690346
(Iteration 3251 / 4900) loss: 0.945628
(Iteration 3301 / 4900) loss: 0.675380
(Iteration 3351 / 4900) loss: 0.730800
(Iteration 3401 / 4900) loss: 0.690676
(Epoch 7 / 10) train acc: 0.747000; val_acc: 0.557000
(Iteration 3451 / 4900) loss: 0.889883
(Iteration 3501 / 4900) loss: 0.670590
(Iteration 3551 / 4900) loss: 0.722332
(Iteration 3601 / 4900) loss: 0.801483
(Iteration 3651 / 4900) loss: 0.670840
(Iteration 3701 / 4900) loss: 0.699421
(Iteration 3751 / 4900) loss: 0.600081
(Iteration 3801 / 4900) loss: 0.816653
(Iteration 3851 / 4900) loss: 0.713218
(Iteration 3901 / 4900) loss: 0.740356
(Epoch 8 / 10) train acc: 0.765000; val_acc: 0.555000
(Iteration 3951 / 4900) loss: 0.623473
(Iteration 4001 / 4900) loss: 0.758548
(Iteration 4051 / 4900) loss: 0.612299
(Iteration 4101 / 4900) loss: 0.677947
(Iteration 4151 / 4900) loss: 0.743771
(Iteration 4201 / 4900) loss: 0.647976
(Iteration 4251 / 4900) loss: 0.664230
(Iteration 4301 / 4900) loss: 0.660933
(Iteration 4351 / 4900) loss: 0.447550
(Iteration 4401 / 4900) loss: 0.688373
(Epoch 9 / 10) train acc: 0.782000; val_acc: 0.546000
(Iteration 4451 / 4900) loss: 0.668955
(Iteration 4501 / 4900) loss: 0.620885
(Iteration 4551 / 4900) loss: 0.647076
(Iteration 4601 / 4900) loss: 0.603572
(Iteration 4651 / 4900) loss: 0.722565
(Iteration 4701 / 4900) loss: 0.511939
(Iteration 4751 / 4900) loss: 0.483064
(Iteration 4801 / 4900) loss: 0.567460
(Iteration 4851 / 4900) loss: 0.557668
(Epoch 10 / 10) train acc: 0.808000; val_acc: 0.558000

```

```

[40]: 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: 0.571

Test set accuracy: 0.559

2 optim.py

```
[ ]: import numpy as np

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu). It has been modified in various areas for use in the
ECE 239AS class at UCLA. This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
permission to use this code. To see the original version, please visit
cs231n.stanford.edu.
"""

"""
This file implements various first-order update rules that are commonly used for
training neural networks. Each update rule accepts current weights and the
gradient of the loss with respect to those weights and produces the next set of
weights. Each update rule has the same interface:

def update(w, dw, config=None):

Inputs:
- w: A numpy array giving the current weights.
- dw: A numpy array of the same shape as w giving the gradient of the
    loss with respect to w.
- config: A dictionary containing hyperparameter values such as learning rate,
    momentum, etc. If the update rule requires caching values over many
    iterations, then config will also hold these cached values.

Returns:
- next_w: The next point after the update.
- config: The config dictionary to be passed to the next iteration of the
    update rule.

NOTE: For most update rules, the default learning rate will probably not perform
well; however the default values of the other hyperparameters should work well
for a variety of different problems.
```

For efficiency, update rules may perform in-place updates, mutating w and setting $next_w$ equal to w .

```
"""
def sgd(w, dw, config=None):
    """
    Performs vanilla stochastic gradient descent.

    config format:
    - learning_rate: Scalar learning rate.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)

    w -= config['learning_rate'] * dw
    return w, config

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 <- alpha*v - epsilon*g
    v = (config['momentum'] * v) - (config['learning_rate'] * dw)
    next_w = w + v

    # ===== #
```

```

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

config['velocity'] = v

return next_w, config

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_new <- alpha*v_old - epsilon*g
    # theta = theta + v_new + alpha*(v_new - v_old)
    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

    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
```

```
def adam(w, dw, config=None):
```

"""

Uses the Adam update rule, which incorporates moving averages of both the gradient and its square and a bias correction term.

config format:

- learning_rate: Scalar learning rate.
- beta1: Decay rate for moving average of first moment of gradient.
- beta2: Decay rate for moving average of second moment of gradient.
- epsilon: Small scalar used for smoothing to avoid dividing by zero.

```

- m: Moving average of gradient.
- v: Moving average of squared gradient.
- t: Iteration number.
"""
if config is None: config = {}
config.setdefault('learning_rate', 1e-3)
config.setdefault('beta1', 0.9)
config.setdefault('beta2', 0.999)
config.setdefault('epsilon', 1e-8)
config.setdefault('v', np.zeros_like(w))
config.setdefault('a', np.zeros_like(w))
config.setdefault('t', 0)

next_w = None

# ===== #
# YOUR CODE HERE:
# Implement Adam. 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, and in config['v'] the moving average of the
# first moments. Finally, store in config['t'] the increasing time.
# ===== #

config['t'] += 1
config['v'] = config['beta1'] * config['v'] + (1 - config['beta1']) * dw
config['a'] = config['beta2'] * config['a'] + (1 - config['beta2']) * np.
→multiply(dw, dw)

v_correct = config['v'] / (1 - config['beta1']**config['t'])
a_correct = config['a'] / (1 - config['beta2']**config['t'])

next_w = w - np.multiply((config['learning_rate'] / (np.sqrt(a_correct) +
→config['epsilon']))), v_correct)

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

return next_w, config

```