

Convolutional neural networks

In this notebook, we'll put together our convolutional layers to implement a 3-layer CNN. Then, we'll ask you to implement a CNN that can achieve > 65% validation error on CIFAR-10.

If you have not completed the Spatial BatchNorm Notebook, please see the following description from that notebook:

Please copy and paste your prior implemented code from HW #4 to start this assignment. If you did not correctly implement the layers in HW #4, you may collaborate with a classmate to use their layer implementations from HW #4. You may also visit TA or Prof OH to correct your implementation.

You'll want to copy and paste from HW #4: - layers.py for your FC network layers, as well as batchnorm and dropout. - layer_utils.py for your combined FC network layers. - optim.py for your optimizers.

Be sure to place these in the `nndl/` directory so they're imported correctly. Note, as announced in class, we will not be releasing our solutions.

```
In [2]: # As usual, a bit of setup

import numpy as np
import matplotlib.pyplot as plt
from nndl.cnn import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient_array, eval_numerical
from nndl.layers import *
from nndl.conv_layers import *
from utils.fast_layers import *
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))))
```

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

```
%reload_ext autoreload
```

```
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,)
```

Three layer CNN

In this notebook, you will implement a three layer CNN. The `ThreeLayerConvNet` class is in `nndl/cnn.py`. You'll need to modify that code for this section, including the initialization, as well as the calculation of the loss and gradients. You should be able to use the building blocks you have either earlier coded or that we have provided. Be sure to use the fast layers.

The architecture of this CNN will be:

conv - relu - 2x2 max pool - affine - relu - affine - softmax

We won't use batchnorm yet. You've also done enough of these to know how to debug; use the cells below.

Note: As we are implementing several layers CNN networks. The gradient error can be expected for the `eval_numerical_gradient()` function. If your `W1 max relative error` and `W2 max relative error` are around or below 0.01, they should be acceptable. Other errors should be less than 1e-5.

```
In [57]: num_inputs = 2
input_dim = (3, 16, 16)
reg = 0.0
num_classes = 10
X = np.random.randn(num_inputs, *input_dim)
y = np.random.randint(num_classes, size=num_inputs)

model = ThreeLayerConvNet(num_filters=3, filter_size=3,
                           input_dim=input_dim, hidden_dim=7,
                           dtype=np.float64)
loss, grads = model.loss(X, y)
for param_name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    param_grad_num = eval_numerical_gradient(f, model.params[param_name], verbose=False)
    e = rel_error(param_grad_num, grads[param_name])
    print('{} max relative error: {}'.format(param_name, rel_error(param_grad_num, grads[param_name])))

W1 max relative error: 0.0008542047071485385
W2 max relative error: 0.0004890474152414171
W3 max relative error: 6.766600429490553e-06
b1 max relative error: 2.217326003638511e-05
b2 max relative error: 3.94244512737717e-08
b3 max relative error: 1.386475678644343e-09
```

Overfit small dataset

To check your CNN implementation, let's overfit a small dataset.

```
In [58]: num_train = 100
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'],
}

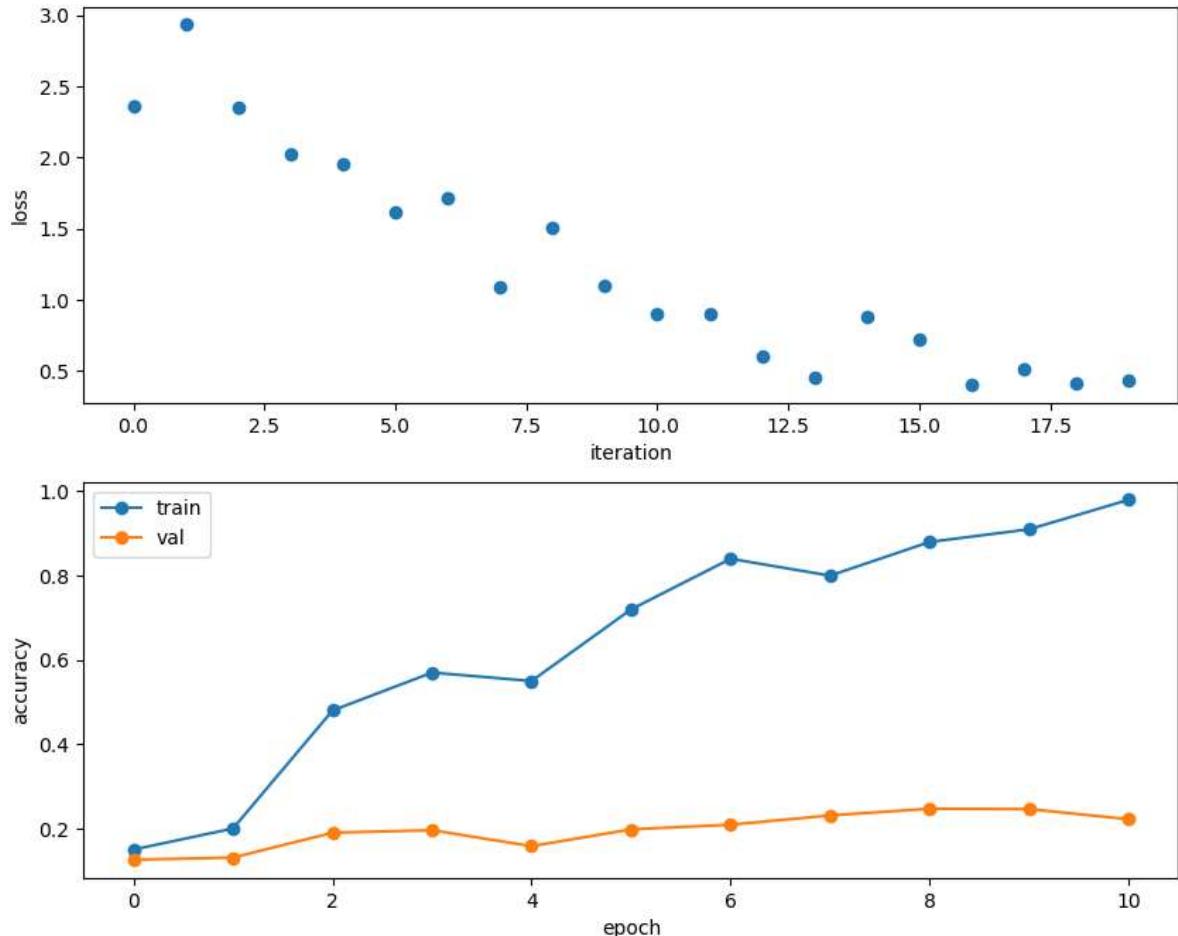
model = ThreeLayerConvNet(weight_scale=1e-2)

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

(Iteration 1 / 20) loss: 2.357512
(Epoch 0 / 10) train acc: 0.150000; val_acc: 0.126000
(Iteration 2 / 20) loss: 2.932379
(Epoch 1 / 10) train acc: 0.200000; val_acc: 0.131000
(Iteration 3 / 20) loss: 2.351277
(Iteration 4 / 20) loss: 2.021465
(Epoch 2 / 10) train acc: 0.480000; val_acc: 0.190000
(Iteration 5 / 20) loss: 1.953078
(Iteration 6 / 20) loss: 1.612692
(Epoch 3 / 10) train acc: 0.570000; val_acc: 0.196000
(Iteration 7 / 20) loss: 1.711191
(Iteration 8 / 20) loss: 1.085145
(Epoch 4 / 10) train acc: 0.550000; val_acc: 0.158000
(Iteration 9 / 20) loss: 1.508036
(Iteration 10 / 20) loss: 1.099567
(Epoch 5 / 10) train acc: 0.720000; val_acc: 0.198000
(Iteration 11 / 20) loss: 0.898851
(Iteration 12 / 20) loss: 0.895582
(Epoch 6 / 10) train acc: 0.840000; val_acc: 0.209000
(Iteration 13 / 20) loss: 0.598804
(Iteration 14 / 20) loss: 0.448797
(Epoch 7 / 10) train acc: 0.800000; val_acc: 0.231000
(Iteration 15 / 20) loss: 0.876180
(Iteration 16 / 20) loss: 0.722372
(Epoch 8 / 10) train acc: 0.880000; val_acc: 0.247000
(Iteration 17 / 20) loss: 0.402560
(Iteration 18 / 20) loss: 0.509233
(Epoch 9 / 10) train acc: 0.910000; val_acc: 0.246000
(Iteration 19 / 20) loss: 0.411652
(Iteration 20 / 20) loss: 0.429783
(Epoch 10 / 10) train acc: 0.980000; val_acc: 0.222000
```

```
In [59]: plt.subplot(2, 1, 1)
plt.plot(solver.loss_history, 'o')
plt.xlabel('iteration')
plt.ylabel('loss')

plt.subplot(2, 1, 2)
plt.plot(solver.train_acc_history, '-o')
plt.plot(solver.val_acc_history, '-o')
plt.legend(['train', 'val'], loc='upper left')
plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.show()
```



Train the network

Now we train the 3 layer CNN on CIFAR-10 and assess its accuracy.

```
In [64]: model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)

solver = Solver(model, data,
                num_epochs=1, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=20)
solver.train()
```

```
(Iteration 1 / 980) loss: 2.304756
(Epoch 0 / 1) train acc: 0.107000; val_acc: 0.112000
(Iteration 21 / 980) loss: 2.202013
(Iteration 41 / 980) loss: 2.049808
(Iteration 61 / 980) loss: 1.878591
(Iteration 81 / 980) loss: 2.100942
(Iteration 101 / 980) loss: 1.763225
(Iteration 121 / 980) loss: 1.923981
(Iteration 141 / 980) loss: 2.048874
(Iteration 161 / 980) loss: 1.809491
(Iteration 181 / 980) loss: 2.014150
(Iteration 201 / 980) loss: 1.724650
(Iteration 221 / 980) loss: 1.710246
(Iteration 241 / 980) loss: 1.636040
(Iteration 261 / 980) loss: 2.084851
(Iteration 281 / 980) loss: 1.930749
(Iteration 301 / 980) loss: 1.771768
(Iteration 321 / 980) loss: 1.667101
(Iteration 341 / 980) loss: 1.646237
(Iteration 361 / 980) loss: 1.772351
(Iteration 381 / 980) loss: 1.936586
(Iteration 401 / 980) loss: 1.374713
(Iteration 421 / 980) loss: 1.646301
(Iteration 441 / 980) loss: 1.573109
(Iteration 461 / 980) loss: 1.712331
(Iteration 481 / 980) loss: 1.435039
(Iteration 501 / 980) loss: 1.504138
(Iteration 521 / 980) loss: 1.742027
(Iteration 541 / 980) loss: 1.421509
(Iteration 561 / 980) loss: 1.536095
(Iteration 581 / 980) loss: 1.587538
(Iteration 601 / 980) loss: 1.735481
(Iteration 621 / 980) loss: 1.810468
(Iteration 641 / 980) loss: 1.674482
(Iteration 661 / 980) loss: 1.714118
(Iteration 681 / 980) loss: 1.588084
(Iteration 701 / 980) loss: 1.735320
(Iteration 721 / 980) loss: 1.575969
(Iteration 741 / 980) loss: 1.240698
(Iteration 761 / 980) loss: 1.409613
(Iteration 781 / 980) loss: 1.418835
(Iteration 801 / 980) loss: 1.537635
(Iteration 821 / 980) loss: 1.304083
(Iteration 841 / 980) loss: 1.664136
(Iteration 861 / 980) loss: 1.502191
(Iteration 881 / 980) loss: 1.632431
(Iteration 901 / 980) loss: 1.469633
(Iteration 921 / 980) loss: 1.744118
(Iteration 941 / 980) loss: 1.402994
(Iteration 961 / 980) loss: 1.669644
(Epoch 1 / 1) train acc: 0.465000; val_acc: 0.469000
```

Get > 65% validation accuracy on CIFAR-10.

In the last part of the assignment, we'll now ask you to train a CNN to get better than 65% validation accuracy on CIFAR-10.

Things you should try:

- Filter size: Above we used 7x7; but VGGNet and onwards showed stacks of 3x3 filters are good.
- Number of filters: Above we used 32 filters. Do more or fewer do better?
- Batch normalization: Try adding spatial batch normalization after convolution layers and vanilla batch normalization after affine layers. Do your networks train faster?
- Network architecture: Can a deeper CNN do better? Consider these architectures:
 - [conv-relu-pool]xN - conv - relu - [affine]xM - [softmax or SVM]
 - [conv-relu-pool]XN - [affine]XM - [softmax or SVM]
 - [conv-relu-conv-relu-pool]xN - [affine]xM - [softmax or SVM]

Tips for training

For each network architecture that you try, you should tune the learning rate and regularization strength. When doing this there are a couple important things to keep in mind:

- If the parameters are working well, you should see improvement within a few hundred iterations
- Remember the coarse-to-fine approach for hyperparameter tuning: start by testing a large range of hyperparameters for just a few training iterations to find the combinations of parameters that are working at all.
- Once you have found some sets of parameters that seem to work, search more finely around these parameters. You may need to train for more epochs.

In [66]:

```
# ===== #
# YOUR CODE HERE:
#   Implement a CNN to achieve greater than 65% validation accuracy
#   on CIFAR-10.
# ===== #
model = ThreeLayerConvNet(weight_scale=0.001,
                           hidden_dim=500,
                           reg=0.001,
                           filter_size=3,
                           use_batchnorm=True)

solver = Solver(model, data,
                num_epochs=10, batch_size=500,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                lr_decay=0.95,
                verbose=True, print_every=20)
solver.train()

y_val_max = np.argmax(model.loss(data['X_val']), axis=1)
print('Validation set accuracy: {}'.format(np.mean(y_val_max == data['y_val'])))
# ===== #
# END YOUR CODE HERE
# ===== #
```

```
(Iteration 1 / 980) loss: 2.304614
(Epoch 0 / 10) train acc: 0.228000; val_acc: 0.224000
(Iteration 21 / 980) loss: 1.596177
(Iteration 41 / 980) loss: 1.426527
(Iteration 61 / 980) loss: 1.435928
(Iteration 81 / 980) loss: 1.287795
(Epoch 1 / 10) train acc: 0.638000; val_acc: 0.566000
(Iteration 101 / 980) loss: 1.246876
(Iteration 121 / 980) loss: 1.181536
(Iteration 141 / 980) loss: 1.147989
(Iteration 161 / 980) loss: 1.137388
(Iteration 181 / 980) loss: 1.094934
(Epoch 2 / 10) train acc: 0.661000; val_acc: 0.625000
(Iteration 201 / 980) loss: 1.064896
(Iteration 221 / 980) loss: 1.032609
(Iteration 241 / 980) loss: 1.014800
(Iteration 261 / 980) loss: 0.921793
(Iteration 281 / 980) loss: 0.956561
(Epoch 3 / 10) train acc: 0.737000; val_acc: 0.650000
(Iteration 301 / 980) loss: 0.929511
(Iteration 321 / 980) loss: 0.885943
(Iteration 341 / 980) loss: 0.962241
(Iteration 361 / 980) loss: 0.845772
(Iteration 381 / 980) loss: 0.886433
(Epoch 4 / 10) train acc: 0.762000; val_acc: 0.657000
(Iteration 401 / 980) loss: 0.837336
(Iteration 421 / 980) loss: 0.892782
(Iteration 441 / 980) loss: 0.832856
(Iteration 461 / 980) loss: 0.815616
(Iteration 481 / 980) loss: 0.835599
(Epoch 5 / 10) train acc: 0.793000; val_acc: 0.668000
(Iteration 501 / 980) loss: 0.734404
(Iteration 521 / 980) loss: 0.838730
(Iteration 541 / 980) loss: 0.734107
(Iteration 561 / 980) loss: 0.715803
(Iteration 581 / 980) loss: 0.765392
(Epoch 6 / 10) train acc: 0.829000; val_acc: 0.685000
(Iteration 601 / 980) loss: 0.707300
(Iteration 621 / 980) loss: 0.692557
(Iteration 641 / 980) loss: 0.654190
(Iteration 661 / 980) loss: 0.653373
(Iteration 681 / 980) loss: 0.650290
(Epoch 7 / 10) train acc: 0.850000; val_acc: 0.699000
(Iteration 701 / 980) loss: 0.700956
(Iteration 721 / 980) loss: 0.674887
(Iteration 741 / 980) loss: 0.601462
(Iteration 761 / 980) loss: 0.643344
(Iteration 781 / 980) loss: 0.648743
(Epoch 8 / 10) train acc: 0.880000; val_acc: 0.665000
(Iteration 801 / 980) loss: 0.598982
(Iteration 821 / 980) loss: 0.618375
(Iteration 841 / 980) loss: 0.624425
(Iteration 861 / 980) loss: 0.612457
(Iteration 881 / 980) loss: 0.607874
(Epoch 9 / 10) train acc: 0.873000; val_acc: 0.681000
(Iteration 901 / 980) loss: 0.546300
(Iteration 921 / 980) loss: 0.499625
```

```
(Iteration 941 / 980) loss: 0.568108
(Iteration 961 / 980) loss: 0.567807
(Epoch 10 / 10) train acc: 0.905000; val_acc: 0.692000
Validation set accuracy: 0.702
```



```
In [ ]: import numpy as np
from nn1.layers import *
import pdb

def conv_forward_naive(x, w, b, conv_param):
    """
    A naive implementation of the forward pass for a convolutional layer.

    The input consists of N data points, each with C channels, height H and width
    W. We convolve each input with F different filters, where each filter spans
    all C channels and has height HH and width WW.

    Input:
    - x: Input data of shape (N, C, H, W)
    - w: Filter weights of shape (F, C, HH, WW)
    - b: Biases, of shape (F,)
    - conv_param: A dictionary with the following keys:
        - 'stride': The number of pixels between adjacent receptive fields in the
                    horizontal and vertical directions.
        - 'pad': The number of pixels that will be used to zero-pad the input.

    Returns a tuple of:
    - out: Output data, of shape (N, F, H', W') where H' and W' are given by
        H' = 1 + (H + 2 * pad - HH) / stride
        W' = 1 + (W + 2 * pad - WW) / stride
    - cache: (x, w, b, conv_param)
    """
    out = None
    pad = conv_param['pad']
    stride = conv_param['stride']

    # ===== #
    # YOUR CODE HERE:
    # Implement the forward pass of a convolutional neural network.
    # Store the output as 'out'.
    # Hint: to pad the array, you can use the function np.pad.
    # ===== #
    N, C, H, W = x.shape
    F, C, HH, WW = w.shape

    H_prime = 1 + (H + 2 * pad - HH) / stride
    W_prime = 1 + (W + 2 * pad - WW) / stride
    H_prime, W_prime = np.int(H_prime), np.int(W_prime)

    xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad, pad)))
    out = np.zeros((N, F, H_prime, W_prime))

    for img in range(N):
        for filt in range(F):
            weights = w[filt]
            for r in range(H_prime):
                for c in range(W_prime):
                    stride_r, stride_c = r*stride, c*stride
                    out[img, filt, r, c] = np.sum(weights * xpad[img, :, stride_r : stride_r + HH, stride_c : stride_c + WW])

    # ===== #

```

```

# END YOUR CODE HERE
# ===== #
cache = (x, w, b, conv_param)
return out, cache

def conv_backward_naive(dout, cache):
    """
    A naive implementation of the backward pass for a convolutional layer.

    Inputs:
    - dout: Upstream derivatives.
    - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive

    Returns a tuple of:
    - dx: Gradient with respect to x
    - dw: Gradient with respect to w
    - db: Gradient with respect to b
    """
dx, dw, db = None, None, None

N, F, out_height, out_width = dout.shape
x, w, b, conv_param = cache

stride, pad = [conv_param['stride'], conv_param['pad']]
xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad,pad)), mode='constant')
num_filts, _, f_height, f_width = w.shape

# ===== #
# YOUR CODE HERE:
# Implement the backward pass of a convolutional neural network.
# Calculate the gradients: dx, dw, and db.
# ===== #
H, W = x.shape[2], x.shape[3]

db = np.zeros(b.shape)
for filt in range(num_filts):
    db[filt] = np.sum(dout[:, filt, ...])

dw = np.zeros(w.shape)
for img in range(N):
    for filt in range(num_filts):
        for r in range(out_height):
            for c in range(out_width):
                stride_r, stride_c = r*stride, c*stride
                dw[filt] += dout[img, filt, r, c] * xpad[img, :, stride_r : stride_r+f_height, stride_c : stride_c+f_width]

dx_nopad = np.zeros(xpad.shape)
for img in range(N):
    for filt in range(num_filts):
        weights = w[filt]
        for r in range(out_height):
            for c in range(out_width):
                stride_r, stride_c = r*stride, c*stride
                dx_nopad[img, :, stride_r : stride_r+f_height, stride_c : stride_c+f_width] += weights[r, c] * dout[img, filt, r, c]

```

```

dx = dx_nopad[..., pad : H+pad, pad : W+pad]
# ===== #
# END YOUR CODE HERE
# ===== #

return dx, dw, db

```



```

def max_pool_forward_naive(x, pool_param):
    """
    A naive implementation of the forward pass for a max pooling layer.

    Inputs:
    - x: Input data, of shape (N, C, H, W)
    - pool_param: dictionary with the following keys:
        - 'pool_height': The height of each pooling region
        - 'pool_width': The width of each pooling region
        - 'stride': The distance between adjacent pooling regions

    Returns a tuple of:
    - out: Output data
    - cache: (x, pool_param)
    """
    out = None

    # ===== #
    # YOUR CODE HERE:
    # Implement the max pooling forward pass.
    # ===== #
    pool_height, pool_width, stride = pool_param['pool_height'], pool_param['pool_width'], pool_param['stride']

    N, C, H, W = x.shape
    out_height = (H - pool_height) / stride + 1
    out_width = (W - pool_width) / stride + 1
    out_height, out_width = np.int(out_height), np.int(out_width)
    out = np.zeros((N, C, out_height, out_width))

    for img in range(N):
        for channel in range(C):
            for r in range(out_height):
                for c in range(out_width):
                    stride_r, stride_c = r*stride, c*stride
                    out[img, channel, r, c] = np.max(x[img, channel, stride_r : stride_r + pool_height, stride_c : stride_c + pool_width])

    # ===== #
    # END YOUR CODE HERE
    # ===== #
    cache = (x, pool_param)
    return out, cache

```



```

def max_pool_backward_naive(dout, cache):
    """
    A naive implementation of the backward pass for a max pooling layer.

    Inputs:
    - dout: Upstream derivatives
    - cache: A tuple of (x, pool_param) as in the forward pass.
    """

```

```

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

dx = None
x, pool_param = cache
pool_height, pool_width, stride = pool_param['pool_height'], pool_param['pool_width']

# ===== #
# YOUR CODE HERE:
#   Implement the max pooling backward pass.
# ===== #
N, C, H, W = x.shape
_, _, dout_height, dout_width = dout.shape
dx = np.zeros(x.shape)

for img in range(N):
    for channel in range(C):
        for r in range(dout_height):
            for c in range(dout_width):
                stride_r, stride_c = r*stride, c*stride
                window = x[img, channel, stride_r : stride_r + pool_height, stride_c : stride_c + pool_width]
                d = dout[img, channel, r, c] * (window == np.max(window))
                dx[img, channel, stride_r : stride_r + pool_height, stride_c : stride_c + pool_width] += d

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

return dx

def spatial_batchnorm_forward(x, gamma, beta, bn_param):
"""
Computes the forward pass for spatial batch normalization.

Inputs:
- x: Input data of shape (N, C, H, W)
- gamma: Scale parameter, of shape (C,)
- beta: Shift parameter, of shape (C,)
- bn_param: Dictionary with the following keys:
    - mode: 'train' or 'test'; required
    - eps: Constant for numeric stability
    - momentum: Constant for running mean / variance. momentum=0 means that
      old information is discarded completely at every time step, while
      momentum=1 means that new information is never incorporated. The
      default of momentum=0.9 should work well in most situations.
    - running_mean: Array of shape (D,) giving running mean of features
    - running_var Array of shape (D,) giving running variance of features

Returns a tuple of:
- out: Output data, of shape (N, C, H, W)
- cache: Values needed for the backward pass
"""
out, cache = None, None

# ===== #
# YOUR CODE HERE:
#   Implement the spatial batchnorm forward pass.
#

```

```

# You may find it useful to use the batchnorm forward pass you
# implemented in HW #4.
# ===== #

N, C, H, W = x.shape
x_transpose = x.transpose(0, 2, 3, 1)
x_reshape = x_transpose.reshape(N*H*W, C)
out, cache = batchnorm_forward(x_reshape, gamma, beta, bn_param)
out = out.reshape(N, H, W, C).transpose(0, 3, 1, 2)
# ===== #
# END YOUR CODE HERE
# ===== #

return out, cache

def spatial_batchnorm_backward(dout, cache):
    """
    Computes the backward pass for spatial batch normalization.

    Inputs:
    - dout: Upstream derivatives, of shape (N, C, H, W)
    - cache: Values from the forward pass

    Returns a tuple of:
    - dx: Gradient with respect to inputs, of shape (N, C, H, W)
    - dgamma: Gradient with respect to scale parameter, of shape (C,)
    - dbeta: Gradient with respect to shift parameter, of shape (C,)
    """
    dx, dgamma, dbeta = None, None, None

    # ===== #
    # YOUR CODE HERE:
    # Implement the spatial batchnorm backward pass.
    #
    # You may find it useful to use the batchnorm forward pass you
    # implemented in HW #4.
    # ===== #

    N, C, H, W = dout.shape
    dout_transpose = dout.transpose(0, 2, 3, 1)
    dout_reshape = dout_transpose.reshape(N*H*W, C)
    dx, dgamma, dbeta = batchnorm_backward(dout_reshape, cache)
    dx = dx.reshape(N, H, W, C).transpose(0, 3, 1, 2)

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

    return dx, dgamma, dbeta

```



```
In [ ]: import numpy as np

from nndl.layers import *
from nndl.conv_layers import *
from utils.fast_layers import *
from nndl.layer_utils import *
from nndl.conv_layer_utils import *

import pdb

class ThreeLayerConvNet(object):
    """
    A three-layer convolutional network with the following architecture:

    conv - relu - 2x2 max pool - affine - relu - affine - softmax

    The network operates on minibatches of data that have shape (N, C, H, W)
    consisting of N images, each with height H and width W and with C input
    channels.
    """

    def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
                 hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
                 dtype=np.float32, use_batchnorm=False):
        """
        Initialize a new network.

        Inputs:
        - input_dim: Tuple (C, H, W) giving size of input data
        - num_filters: Number of filters to use in the convolutional layer
        - filter_size: Size of filters to use in the convolutional layer
        - hidden_dim: Number of units to use in the fully-connected hidden layer
        - num_classes: Number of scores to produce from the final affine layer.
        - weight_scale: Scalar giving standard deviation for random initialization
          of weights.
        - reg: Scalar giving L2 regularization strength
        - dtype: numpy datatype to use for computation.
        """
        self.use_batchnorm = use_batchnorm
        self.params = {}
        self.reg = reg
        self.dtype = dtype

        # ===== #
        # YOUR CODE HERE:
        #   Initialize the weights and biases of a three Layer CNN. To initialize:
        #       - the biases should be initialized to zeros.
        #       - the weights should be initialized to a matrix with entries
        #           drawn from a Gaussian distribution with zero mean and
        #           standard deviation given by weight_scale.
        # ===== #

        C, H, W = input_dim
        stride = 1
        pad = (filter_size - 1) / 2
```

```

output_conv_height = (H + 2*pad - filter_size) / stride + 1
output_conv_width = (W + 2*pad - filter_size) / stride + 1

self.params['W1'] = np.random.normal(0, weight_scale, (num_filters, C, fil
self.params['b1'] = np.zeros(num_filters)

# 2x2 max pool, stride 2
output_pool_height = int((output_conv_height - 2) / 2 + 1)
output_pool_width = int((output_conv_height - 2) / 2 + 1)

pool_size = output_pool_height * output_pool_width * num_filters
self.params['W2'] = np.random.normal(0, weight_scale, (pool_size, hidden_d
self.params['b2'] = np.zeros(hidden_dim)

self.params['W3'] = np.random.normal(0, weight_scale, (hidden_dim, num_cla
self.params['b3'] = np.zeros(num_classes)

self.bn_params = {}
if self.use_batchnorm:
    self.bn_params['bn_param1'] = {'mode': 'train'}
    self.params['beta1'] = np.zeros(num_filters)
    self.params['gamma1'] = np.ones(num_filters)

    self.bn_params['bn_param2'] = {'mode': 'train'}
    self.params['beta2'] = np.zeros(hidden_dim)
    self.params['gamma2'] = np.ones(hidden_dim)

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

for k, v in self.params.items():
    self.params[k] = v.astype(dtype)

def loss(self, X, y=None):
    """
    Evaluate loss and gradient for the three-layer convolutional network.

    Input / output: Same API as TwoLayerNet in fc_net.py.
    """
    W1, b1 = self.params['W1'], self.params['b1']
    W2, b2 = self.params['W2'], self.params['b2']
    W3, b3 = self.params['W3'], self.params['b3']

    # pass conv_param to the forward pass for the convolutional layer
    filter_size = W1.shape[2]
    conv_param = {'stride': 1, 'pad': (filter_size - 1) / 2}

    # pass pool_param to the forward pass for the max-pooling Layer
    pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

    scores = None

    # ===== #
    # YOUR CODE HERE:
    # Implement the forward pass of the three Layer CNN. Store the output

```

```

#      scores as the variable "scores".
# ===== #
if self.use_batchnorm:
    bn_param1 = self.bn_params['bn_param1']
    bn_param2 = self.bn_params['bn_param2']

    beta1 = self.params['beta1']
    beta2 = self.params['beta2']

    gamma1 = self.params['gamma1']
    gamma2 = self.params['gamma2']

    c_out, c_cache = conv_batchnorm_relu_pool_forward(X, W1, b1, conv_param,
                                                       N, F, H, W = c_out.shape

    c_out = c_out.reshape((N, F*H*W))
    ar_out, ar_cache = affine_batchnorm_relu_forward(c_out, W2, b2, gamma2)
else:
    c_out, c_cache = conv_relu_pool_forward(X, W1, b1, conv_param, pool_param
                                             ar_out, ar_cache = affine_relu_forward(c_out, W2, b2)

scores, scores_cache = affine_forward(ar_out, W3, b3)

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

if y is None:
    return scores

loss, grads = 0, {}
# ===== #
# YOUR CODE HERE:
# Implement the backward pass of the three Layer CNN. Store the grads
# in the grads dictionary, exactly as before (i.e., the gradient of
# self.params[k] will be grads[k]). Store the loss as "Loss", and
# don't forget to add regularization on ALL weight matrices.
# ===== #

loss, grad_loss = softmax_loss(scores, y)

loss += 0.5 * self.reg * (np.sum(np.square(W1)) + np.sum(np.square(W2)) + ...

dx, dw3, grads['b3'] = affine_backward(grad_loss, scores_cache)

if self.use_batchnorm:
    dx, dw2, grads['b2'], grads['gamma2'], grads['beta2'] = affine_batchnorm_backward(dx, ar_cache)
    dx = dx.reshape((N, F, H, W))
    dx, dw1, grads['b1'], grads['gamma1'], grads['beta1'] = conv_batchnorm_backward(dx, c_cache)
else:
    dx, dw2, grads['b2'] = affine_relu_backward(dx, ar_cache)
    dx, dw1, grads['b1'] = conv_relu_pool_backward(dx, c_cache)

grads['W3'] = dw3 + self.reg * W3
grads['W2'] = dw2 + self.reg * W2
grads['W1'] = dw1 + self.reg * W1

```

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

return loss, grads

pass
```

```
In [ ]: def conv_batchnorm_relu_pool_forward(x, w, b, conv_param, pool_param, gamma, beta):
    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
    out, bn_cache = spatial_batchnorm_forward(a, gamma, beta, bn_param)
    s, relu_cache = relu_forward(out)
    out, pool_cache = max_pool_forward_fast(s, pool_param)
    cache = (conv_cache, bn_cache, relu_cache, pool_cache)
    return out, cache

def conv_batchnorm_relu_pool_backward(dout, cache):
    conv_cache, bn_cache, relu_cache, pool_cache = cache
    ds = max_pool_backward_fast(dout, pool_cache)
    da = relu_backward(ds, relu_cache)
    dx, dgamma, dbeta = spatial_batchnorm_backward(da, bn_cache)
    dx, dw, db = conv_backward_fast(dx, conv_cache)
    return dx, dw, db, dgamma, dbeta
```