

# Dropout

October 20, 2019

## 1 Dropout

Dropout [1] is a technique for regularizing neural networks by randomly setting some features to zero during the forward pass. In this exercise you will implement a dropout layer and modify your fully-connected network to optionally use dropout.

[1] [Geoffrey E. Hinton et al, "Improving neural networks by preventing co-adaptation of feature detectors", arXiv 2012](<https://arxiv.org/abs/1207.0580>)

```
In [1]: # As usual, a bit of setup
        from __future__ import print_function
        import time
        import numpy as np
        import matplotlib.pyplot as plt
        from cs682.classifiers.fc_net import *
        from cs682.data_utils import get_CIFAR10_data
        from cs682.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
        from cs682.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, v in data.items():
            print('%s: ' % k, v.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,)
```

## 2 Dropout forward pass

In the file `cs682/layers.py`, implement the forward pass for dropout. Since dropout behaves differently during training and testing, make sure to implement the operation for both modes.

Once you have done so, run the cell below to test your implementation.

```
In [7]: np.random.seed(231)
        x = np.random.randn(500, 500) + 10

        for p in [0.25, 0.4, 0.7]:
            out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
            out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})

            print('Running tests with p = ', p)
            print('Mean of input: ', x.mean())
            print('Mean of train-time output: ', out.mean())
            print('Mean of test-time output: ', out_test.mean())
            print('Fraction of train-time output set to zero: ', (out == 0).mean())
            print('Fraction of test-time output set to zero: ', (out_test == 0).mean())
            print()
```

```
Running tests with p = 0.25
Mean of input: 10.000207878477502
Mean of train-time output: 10.014059116977283
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.749784
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.4
Mean of input: 10.000207878477502
Mean of train-time output: 9.977917658761159
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.600796
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.7
Mean of input: 10.000207878477502
Mean of train-time output: 9.987811912159426
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.30074
Fraction of test-time output set to zero: 0.0
```

## 3 Dropout backward pass

In the file `cs682/layers.py`, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

```
In [8]: np.random.seed(231)
        x = np.random.randn(10, 10) + 10
        dout = np.random.randn(*x.shape)

        dropout_param = {'mode': 'train', 'p': 0.2, 'seed': 123}
        out, cache = dropout_forward(x, dropout_param)
        dx = dropout_backward(dout, cache)
        dx_num = eval_numerical_gradient_array(lambda xx: dropout_forward(xx, dropout_param)[0],
        x, dout)
        # print(dx_num)
```

```
# Error should be around e-10 or less
print('dx relative error: ', rel_error(dx, dx_num))
```

```
dx relative error: 5.44560814873387e-11
```

### 3.1 Inline Question 1:

What happens if we do not divide the values being passed through inverse dropout by  $p$  in the dropout layer? Why does that happen?

### 3.2 Answer:

The outputs at test time will be smaller than they otherwise would be by a factor of  $p$  because we do not drop any neurons at test time. By multiplying by a factor of  $1/p$  during training, we essentially cancel this scaling so the outputs have the same expected value.

## 4 Fully-connected nets with Dropout

In the file `cs682/classifiers/fc_net.py`, modify your implementation to use dropout. Specifically, if the constructor of the net receives a value that is not 1 for the dropout parameter, then the net should add dropout immediately after every ReLU nonlinearity. After doing so, run the following to numerically gradient-check your implementation.

```
In [9]: # np.random.seed(231)
        N, D, H1, H2, C = 2, 15, 20, 30, 10
        X = np.random.randn(N, D)
        y = np.random.randint(C, size=(N,))

        for dropout in [1, 0.75, 0.5]:
            print('Running check with dropout = ', dropout)
            model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                                      weight_scale=5e-2, dtype=np.float64,
                                      dropout=dropout, seed=123)

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

            # Relative errors should be around e-6 or less; Note that it's fine
            # if for dropout=1 you have W2 error be on the order of e-5.
            for name in sorted(grads):
                f = lambda _: model.loss(X, y)[0]
                grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
                print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
            print()
```

```
Running check with dropout = 1
Initial loss: 2.3051948273987857
W1 relative error: 5.25e-07
W2 relative error: 1.89e-05
W3 relative error: 2.92e-07
b1 relative error: 1.34e-07
b2 relative error: 7.09e-08
b3 relative error: 1.49e-10
```

```
Running check with dropout = 0.75
```

```

Initial loss: 2.29898614757146
W1 relative error: 9.74e-07
W2 relative error: 4.12e-08
W3 relative error: 2.89e-08
b1 relative error: 9.62e-09
b2 relative error: 1.90e-09
b3 relative error: 8.93e-11

Running check with dropout = 0.5
Initial loss: 2.302437587710995
W1 relative error: 4.55e-08
W2 relative error: 2.97e-08
W3 relative error: 5.99e-07
b1 relative error: 1.87e-08
b2 relative error: 1.38e-09
b3 relative error: 1.84e-10

```

## 5 Regularization experiment

As an experiment, we will train a pair of two-layer networks on 500 training examples: one will use no dropout, and one will use a keep probability of 0.25. We will then visualize the training and validation accuracies of the two networks over time.

```

In [10]: # Train two identical nets, one with dropout and one without
np.random.seed(231)
num_train = 500
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 = {}
dropout_choices = [1, 0.25]
for dropout in dropout_choices:
    model = FullyConnectedNet([500], dropout=dropout)
    print(dropout)

    solver = Solver(model,
                    small_data,
                    num_epochs=25,
                    batch_size=100,
                    update_rule='adam',
                    optim_config={'learning_rate': 5e-4,
                                },
                    verbose=True,
                    print_every=100)

    solver.train()
    solvers[dropout] = solver

1
(Iteration 1 / 125) loss: 7.856643
(Epoch 0 / 25) train acc: 0.260000; val_acc: 0.184000
(Epoch 1 / 25) train acc: 0.416000; val_acc: 0.258000
(Epoch 2 / 25) train acc: 0.482000; val_acc: 0.276000
(Epoch 3 / 25) train acc: 0.532000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.600000; val_acc: 0.271000
(Epoch 5 / 25) train acc: 0.708000; val_acc: 0.299000

```

```

(Epoch 6 / 25) train acc: 0.722000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.832000; val_acc: 0.255000
(Epoch 8 / 25) train acc: 0.878000; val_acc: 0.269000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.275000
(Epoch 10 / 25) train acc: 0.888000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.926000; val_acc: 0.278000
(Epoch 12 / 25) train acc: 0.960000; val_acc: 0.302000
(Epoch 13 / 25) train acc: 0.964000; val_acc: 0.305000
(Epoch 14 / 25) train acc: 0.966000; val_acc: 0.309000
(Epoch 15 / 25) train acc: 0.976000; val_acc: 0.288000
(Epoch 16 / 25) train acc: 0.988000; val_acc: 0.301000
(Epoch 17 / 25) train acc: 0.988000; val_acc: 0.305000
(Epoch 18 / 25) train acc: 0.990000; val_acc: 0.308000
(Epoch 19 / 25) train acc: 0.988000; val_acc: 0.311000
(Epoch 20 / 25) train acc: 0.990000; val_acc: 0.311000
(Iteration 101 / 125) loss: 0.006070
(Epoch 21 / 25) train acc: 0.998000; val_acc: 0.313000
(Epoch 22 / 25) train acc: 0.976000; val_acc: 0.322000
(Epoch 23 / 25) train acc: 0.986000; val_acc: 0.314000
(Epoch 24 / 25) train acc: 0.990000; val_acc: 0.310000
(Epoch 25 / 25) train acc: 0.994000; val_acc: 0.305000
0.25
(Iteration 1 / 125) loss: 17.318478
(Epoch 0 / 25) train acc: 0.230000; val_acc: 0.177000
(Epoch 1 / 25) train acc: 0.378000; val_acc: 0.243000
(Epoch 2 / 25) train acc: 0.402000; val_acc: 0.254000
(Epoch 3 / 25) train acc: 0.502000; val_acc: 0.276000
(Epoch 4 / 25) train acc: 0.528000; val_acc: 0.298000
(Epoch 5 / 25) train acc: 0.562000; val_acc: 0.297000
(Epoch 6 / 25) train acc: 0.626000; val_acc: 0.290000
(Epoch 7 / 25) train acc: 0.628000; val_acc: 0.298000
(Epoch 8 / 25) train acc: 0.686000; val_acc: 0.310000
(Epoch 9 / 25) train acc: 0.722000; val_acc: 0.289000
(Epoch 10 / 25) train acc: 0.724000; val_acc: 0.300000
(Epoch 11 / 25) train acc: 0.760000; val_acc: 0.305000
(Epoch 12 / 25) train acc: 0.772000; val_acc: 0.280000
(Epoch 13 / 25) train acc: 0.814000; val_acc: 0.303000
(Epoch 14 / 25) train acc: 0.814000; val_acc: 0.341000
(Epoch 15 / 25) train acc: 0.856000; val_acc: 0.352000
(Epoch 16 / 25) train acc: 0.838000; val_acc: 0.303000
(Epoch 17 / 25) train acc: 0.840000; val_acc: 0.291000
(Epoch 18 / 25) train acc: 0.844000; val_acc: 0.315000
(Epoch 19 / 25) train acc: 0.864000; val_acc: 0.325000
(Epoch 20 / 25) train acc: 0.862000; val_acc: 0.308000
(Iteration 101 / 125) loss: 5.259476
(Epoch 21 / 25) train acc: 0.896000; val_acc: 0.320000
(Epoch 22 / 25) train acc: 0.872000; val_acc: 0.298000
(Epoch 23 / 25) train acc: 0.906000; val_acc: 0.317000
(Epoch 24 / 25) train acc: 0.906000; val_acc: 0.320000
(Epoch 25 / 25) train acc: 0.910000; val_acc: 0.325000

```

In [11]: # Plot train and validation accuracies of the two models

```

train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

```

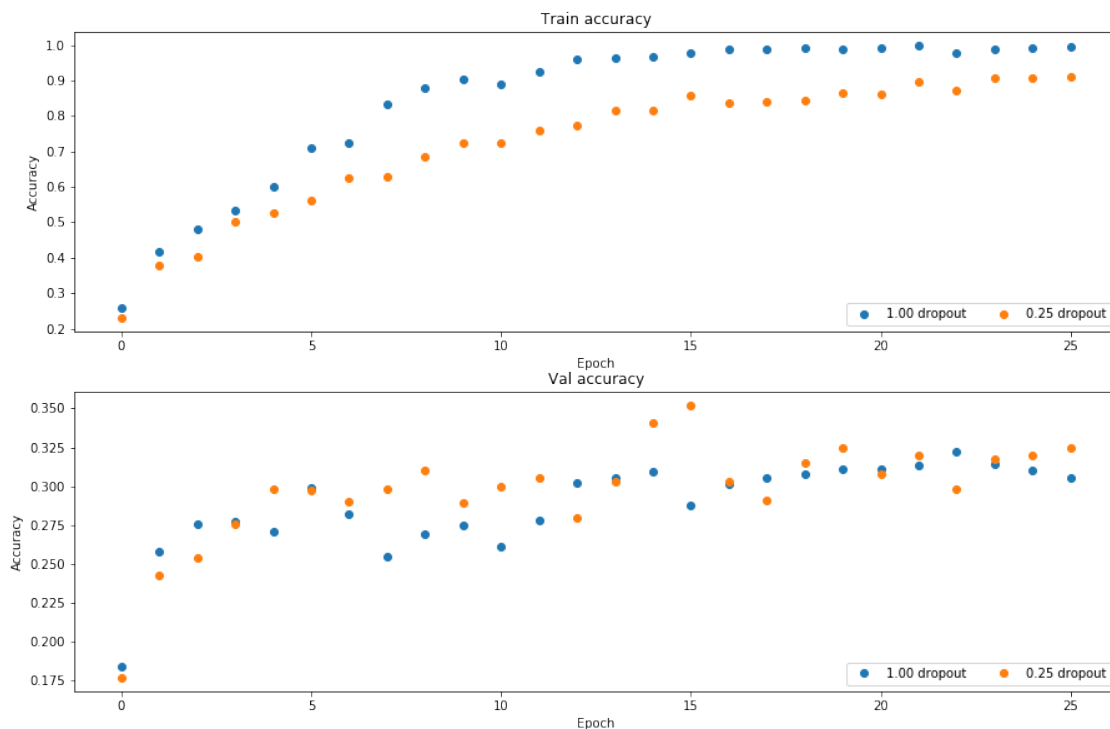
```

plt.subplot(3, 1, 1)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].train_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].val_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Val accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

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

```



## 5.1 Inline Question 2:

Compare the validation and training accuracies with and without dropout -- what do your results suggest about dropout as a regularizer?

## 5.2 Answer:

The network using dropout has a lower training accuracy than the one without, while they have similar validation accuracies, suggesting that dropout is an effective regularizer because it prevents the network from over-fitting on the training data.

### 5.3 Inline Question 3:

Suppose we are training a deep fully-connected network for image classification, with dropout after hidden layers (parameterized by keep probability  $p$ ). How should we modify  $p$ , if at all, if we decide to decrease the size of the hidden layers (that is, the number of nodes in each layer)?

### 5.4 Answer:

We should not modify  $p$ .