

# Softmax exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission. For more details see the [assignments page](http://vision.stanford.edu/teaching/cs231n/assignments.html) (<http://vision.stanford.edu/teaching/cs231n/assignments.html>) on the course website.

This exercise is analogous to the SVM exercise. You will:

- implement a fully-vectorized **loss function** for the Softmax classifier
- implement the fully-vectorized expression for its **analytic gradient**
- **check your implementation** with numerical gradient
- use a validation set to **tune the learning rate and regularization** strength
- **optimize** the loss function with **SGD**
- **visualize** the final learned weights

In [66]:

```
import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

from __future__ import print_function

%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
```

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

```
%reload_ext autoreload
```

In [67]:

```
def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000, num_dev=500):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare it for the linear classifier. These are the same steps as we used for the SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = 'cs231n/datasets/cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # subsample the data
```

```

mask = list(range(num_training, num_training + num_validation))
X_val = X_train[mask]
y_val = y_train[mask]
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]
mask = np.random.choice(num_training, num_dev, replace=False)
X_dev = X_train[mask]
y_dev = y_train[mask]

```

```

# Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

```

```

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

```

```

# add bias dimension and transform into columns
X_train = np.hstack([X_train, np.ones((X_train.shape[0], 1))])
X_val = np.hstack([X_val, np.ones((X_val.shape[0], 1))])
X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
X_dev = np.hstack([X_dev, np.ones((X_dev.shape[0], 1))])

```

```

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

```

```

# Invoke the above function to get our data.

```

```

X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev = get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
print('dev data shape: ', X_dev.shape)
print('dev labels shape: ', y_dev.shape)

```

```

Train data shape: (49000, 3073)
Train labels shape: (49000,)
Validation data shape: (1000, 3073)
Validation labels shape: (1000,)
Test data shape: (1000, 3073)
Test labels shape: (1000,)
dev data shape: (500, 3073)
dev labels shape: (500,)

```

# Softmax Classifier

Your code for this section will all be written inside **cs231n/classifiers/softmax.py**.

In [68]:

```
# First implement the naive softmax loss function with nested loops.
# Open the file cs231n/classifiers/softmax.py and implement the
# softmax_loss_naive function.
from cs231n.classifiers.softmax import softmax_loss_naive
import time

# Generate a random softmax weight matrix and use it to compute the loss.
W = np.random.randn(3073, 10) * 0.0001
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)

# As a rough sanity check, our loss should be something close to -log(0.1).
print('loss: %f' % loss)
print('sanity check: %f' % (-np.log(0.1)))
```

```
loss: 2.407194
sanity check: 2.302585
```

## Inline Question 1:

Why do we expect our loss to be close to  $-\log(0.1)$ ? Explain briefly.\*\*

**Your answer:** A probabilistic interpretation is to look at the expression  $p(y_i|x_i;W) = \frac{\exp(f_{yi})}{\sum_j \exp(f_{ji})}$  as the (normalized) probability assigned to the correct label  $y_i$  given the image  $x_i$  and parameterized by  $W$ . To see this, we remember that the Softmax classifier interprets the scores inside the output vector  $f$  as the unnormalized log probabilities. Exponentiating these quantities therefore gives the (unnormalized) probabilities, and the division performs the normalization so that the probabilities sum to one. In the probabilistic interpretation, we are therefore minimizing the negative log likelihood of the correct class, which can be interpreted as performing Maximum Likelihood Estimation (MLE).  $W$  was chosen randomly, so we have no reason to believe it is more likely to choose the correct label than any other label for any  $x_i$ . Since we have 10 labels we would expect it to be right 10% of the time. So we expect on average to be right 10% of the times. This is exactly the 0.1 inside the  $-\log$ .

In [69]:

```
# Complete the implementation of softmax_loss_naive and implement a (naive)  
# version of the gradient that uses nested loops.  
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 0.0)  
  
# As we did for the SVM, use numeric gradient checking as a debugging tool.  
# The numeric gradient should be close to the analytic gradient.  
from cs231n.gradient_check import grad_check_sparse  
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 0.0)[0]  
grad_numerical = grad_check_sparse(f, W, grad, 10)  
  
# similar to SVM case, do another gradient check with regularization  
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 5e1)  
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 5e1)[0]  
grad_numerical = grad_check_sparse(f, W, grad, 10)
```

numerical: 0.014854 analytic: 0.014854, relative error: 6.806826e-10  
numerical: -1.195550 analytic: -1.195550, relative error: 5.887579e-08  
numerical: 1.906184 analytic: 1.906184, relative error: 5.406220e-08  
numerical: -1.806473 analytic: -1.806473, relative error: 3.218486e-08  
numerical: -0.569003 analytic: -0.569003, relative error: 1.103770e-07  
numerical: -0.117273 analytic: -0.117273, relative error: 1.708623e-07  
numerical: -0.994115 analytic: -0.994115, relative error: 1.812086e-08  
numerical: -0.978387 analytic: -0.978387, relative error: 3.162485e-08  
numerical: -2.286827 analytic: -2.286827, relative error: 1.772868e-08  
numerical: -0.091875 analytic: -0.091875, relative error: 4.225597e-08  
numerical: 5.210453 analytic: 5.210453, relative error: 4.790185e-09  
numerical: -1.258512 analytic: -1.258512, relative error: 2.669814e-08  
numerical: -0.770884 analytic: -0.770884, relative error: 2.144870e-08  
numerical: 1.745232 analytic: 1.745232, relative error: 8.806938e-09  
numerical: 0.103940 analytic: 0.103940, relative error: 3.852408e-07  
numerical: 1.919603 analytic: 1.919603, relative error: 2.759487e-08  
numerical: 1.669794 analytic: 1.669793, relative error: 3.912509e-08  
numerical: 2.418603 analytic: 2.418603, relative error: 1.031937e-08  
numerical: -0.104148 analytic: -0.104148, relative error: 7.306827e-08  
numerical: 3.195089 analytic: 3.195089, relative error: 1.442025e-08

In [70]:

```
# Now that we have a naive implementation of the softmax loss function and its  
gradient,  
# implement a vectorized version in softmax_loss_vectorized.  
# The two versions should compute the same results, but the vectorized version  
should be  
# much faster.  
tic = time.time()  
loss_naive, grad_naive = softmax_loss_naive(W, X_dev, y_dev, 0.000005)  
toc = time.time()  
print('naive loss: %e computed in %fs' % (loss_naive, toc - tic))  
  
from cs231n.classifiers.softmax import softmax_loss_vectorized  
tic = time.time()  
loss_vectorized, grad_vectorized = softmax_loss_vectorized(W, X_dev, y_dev, 0.  
000005)  
toc = time.time()  
print('vectorized loss: %e computed in %fs' % (loss_vectorized, toc - tic))  
  
# As we did for the SVM, we use the Frobenius norm to compare the two versions  
# of the gradient.  
grad_difference = np.linalg.norm(grad_naive - grad_vectorized, ord='fro')  
print('Loss difference: %f' % np.abs(loss_naive - loss_vectorized))  
print('Gradient difference: %f' % grad_difference)
```

```
naive loss: 2.407194e+00 computed in 0.141369s  
vectorized loss: 2.407194e+00 computed in 0.008652s  
Loss difference: 0.000000  
Gradient difference: 0.000000
```

In [75]:

```
# Use the validation set to tune hyperparameters (regularization strength and
# learning rate). You should experiment with different ranges for the learning
# rates and regularization strengths; if you are careful you should be able to
# get a classification accuracy of over 0.35 on the validation set.
from cs231n.classifiers import Softmax
results = {}
best_val = -1
best_softmax = None
learning_rates = [1e-7, 5e-7]
regularization_strengths = [2.5e4, 5e4]

#####
##
# TODO:
#
# Use the validation set to set the learning rate and regularization strength.
#
# This should be identical to the validation that you did for the SVM; save
#
# the best trained softmax classifier in best_softmax.
#
#####
##
for lr in learning_rates:
    for rs in regularization_strengths:
        softmax = Softmax()
        softmax.train(X_train, y_train, learning_rate=lr, reg=rs, num_iters=15
00, verbose=False)
        y_pred = softmax.predict(X_train)
        train_accuracy = np.mean(y_pred==y_train)
        y_val_pred = softmax.predict(X_val)
        val_accuracy = np.mean(y_val_pred==y_val)
        results[(lr,rs)] = train_accuracy, val_accuracy

        if val_accuracy > best_val:
            best_val = val_accuracy
            best_softmax = softmax
#####
##
#
#                                     END OF YOUR CODE
#
#####
##

# Print out results.
for lr, reg in sorted(results):
    train_accuracy, val_accuracy = results[(lr, reg)]
    print('lr %e reg %e train accuracy: %f val accuracy: %f' % (
        lr, reg, train_accuracy, val_accuracy))

print('best validation accuracy achieved during cross-validation: %f' % best_v
al)
```

```
lr 1.000000e-07 reg 2.500000e+04 train accuracy: 0.330469 val accu
racy: 0.348000
lr 1.000000e-07 reg 5.000000e+04 train accuracy: 0.310286 val accu
racy: 0.326000
lr 5.000000e-07 reg 2.500000e+04 train accuracy: 0.327286 val accu
racy: 0.347000
lr 5.000000e-07 reg 5.000000e+04 train accuracy: 0.292653 val accu
racy: 0.309000
best validation accuracy achieved during cross-validation: 0.34800
0
```

In [76]:

```
# evaluate on test set
# Evaluate the best softmax on test set
y_test_pred = best_softmax.predict(X_test)
test_accuracy = np.mean(y_test == y_test_pred)
print('softmax on raw pixels final test set accuracy: %f' % (test_accuracy, ))
```

```
softmax on raw pixels final test set accuracy: 0.339000
```

In [77]:

```
# Visualize the learned weights for each class
w = best_softmax.W[:-1,:] # strip out the bias
w = w.reshape(32, 32, 3, 10)

w_min, w_max = np.min(w), np.max(w)

classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'shi
p', 'truck']
for i in range(10):
    plt.subplot(2, 5, i + 1)

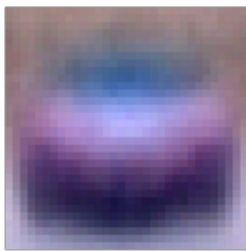
    # Rescale the weights to be between 0 and 255
    wimg = 255.0 * (w[:, :, :, i].squeeze() - w_min) / (w_max - w_min)
    plt.imshow(wimg.astype('uint8'))
    plt.axis('off')
    plt.title(classes[i])
```



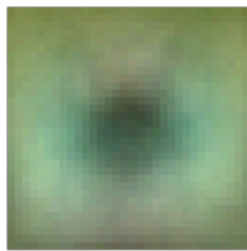
plane



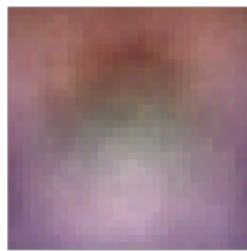
car



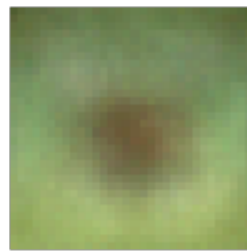
bird



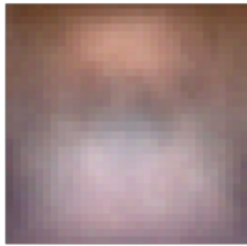
cat



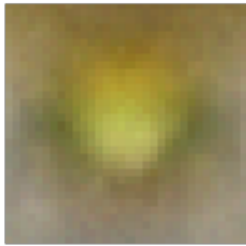
deer



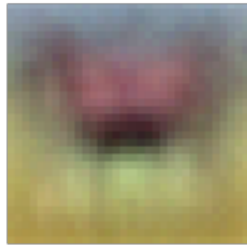
dog



frog



horse



ship



truck

