

# Softmax exercise

Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission.

In this 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 learned weights

Acknowledgement: This exercise is adapted from [Stanford CS231n](#).

```
In [26]: import random
import numpy as np
from data_utils import load_CIFAR10
import matplotlib.pyplot as plt

%matplotlib inline
plt.rcParams['figure.figsize'] = (10, 8) # 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 [27]: def relu_error(out, correct_out):
    return np.sum(abs(out - correct_out)) / (abs(out) + abs(correct_out))

In [28]: 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
    # Softmax, but condensed to a single function.

    # Load the raw CIFAR-10 data
    cifar10_dir = 'datasets/cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # Subsample the data
    mask = range(num_training + num_validation)
    X_val = X_train[mask]
    y_val = y_train[mask]
    mask = range(num_training)
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = range(num_test)
    X_test = X_test[mask]
    y_test = y_test[mask]

    # We now have a development set, which is a small subset of
    # the training data
    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))

    # Compute 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: (500,3)

In [29]: # Create one-hot vectors for label
num_class = 10
y_train_oh = np.zeros((y_train.shape[0], 10))
y_train_oh[np.arange(y_train.shape[0], y_train) == 1] = 1
y_val_oh = np.zeros((y_val.shape[0], 10))
y_val_oh[np.arange(y_val.shape[0], y_val) == 1] = 1
y_test_oh = np.zeros((y_test.shape[0], 10))
y_test_oh[np.arange(y_test.shape[0], y_test) == 1] = 1
y_dev_oh = np.zeros((y_dev.shape[0], 10))
y_dev_oh[np.arange(y_dev.shape[0], y_dev) == 1] = 1
```

## Regression as classifier

The most simple and straightforward approach to learn a classifier is to map the input data (raw image values) to class label (one-hot vector). The loss function is defined as follows:

$$\mathcal{L} = \frac{1}{n} \| \mathbf{XW} - \mathbf{y} \|_2^2 \quad (1)$$

Where:

- $\mathbf{W} \in \mathbb{R}^{(d+1) \times C}$ : Classifier weight
- $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$ : Dataset
- $\mathbf{y} \in \mathbb{R}^{n \times C}$ : Class label (one-hot vector)

## Optimization

Given the loss function (1), the next problem is how to solve the weight  $\mathbf{W}$ . We now discuss 2 approaches:

- Random search
- Closed-form solution

### Random search

```
In [30]: best_loss = float('inf')
for num in range(1001):
    W = np.random.randn(3073, 10) * 0.0001
    loss = np.linalg.norm(X_dev.dot(W) - y_dev_oh)
    if loss < best_loss:
        best_loss = loss
        bestW = W
    print('In attempt %d the loss was %f, best %f' % (num, loss, best_loss))

In attempt 0 the loss was 36.637351, best 36.637351
In attempt 1 the loss was 32.924404, best 32.924404
In attempt 2 the loss was 33.683036, best 32.924404
In attempt 3 the loss was 32.562024, best 32.562024
In attempt 4 the loss was 34.541393, best 32.562024
In attempt 5 the loss was 38.085818, best 32.562024
In attempt 6 the loss was 33.462305, best 32.562024
In attempt 7 the loss was 34.845456, best 32.562024
In attempt 8 the loss was 34.530255, best 32.562024
In attempt 9 the loss was 30.556467, best 30.556467
In attempt 10 the loss was 32.652475, best 30.556467
In attempt 11 the loss was 31.969817, best 30.556467
In attempt 12 the loss was 34.242383, best 30.556467
In attempt 13 the loss was 32.266101, best 30.556467
In attempt 14 the loss was 32.913627, best 30.556467
In attempt 15 the loss was 31.902112, best 30.556467
In attempt 16 the loss was 32.662677, best 30.556467
In attempt 17 the loss was 32.086070, best 30.556467
In attempt 18 the loss was 32.425247, best 30.556467
In attempt 19 the loss was 34.759837, best 30.556467
In attempt 20 the loss was 34.839806, best 30.556467
In attempt 21 the loss was 35.810486, best 30.556467
In attempt 22 the loss was 31.952169, best 30.556467
In attempt 23 the loss was 34.050328, best 30.556467
In attempt 24 the loss was 33.007395, best 30.556467
In attempt 25 the loss was 31.268496, best 30.556467
In attempt 26 the loss was 32.748833, best 30.556467
In attempt 27 the loss was 33.876682, best 30.556467
In attempt 28 the loss was 33.987584, best 30.556467
In attempt 29 the loss was 34.719678, best 30.556467
In attempt 30 the loss was 34.03636, best 30.556467
In attempt 31 the loss was 32.155133, best 30.556467
In attempt 32 the loss was 32.889039, best 30.556467
In attempt 33 the loss was 30.567077, best 30.556467
In attempt 34 the loss was 34.027629, best 30.556467
In attempt 35 the loss was 32.940843, best 30.556467
In attempt 36 the loss was 32.943750, best 30.556467
In attempt 37 the loss was 32.471462, best 30.556467
In attempt 38 the loss was 31.570291, best 30.556467
In attempt 39 the loss was 34.850495, best 30.556467
In attempt 40 the loss was 32.524384, best 30.556467
In attempt 41 the loss was 34.357354, best 30.556467
In attempt 42 the loss was 32.578902, best 30.556467
In attempt 43 the loss was 33.654832, best 30.556467
In attempt 44 the loss was 31.913211, best 30.556467
In attempt 45 the loss was 34.355712, best 30.556467
In attempt 46 the loss was 31.704243, best 30.556467
In attempt 47 the loss was 33.175149, best 30.556467
In attempt 48 the loss was 34.071437, best 30.556467
In attempt 49 the loss was 34.717653, best 30.556467
In attempt 50 the loss was 33.740219, best 30.556467
In attempt 51 the loss was 33.183698, best 30.556467
In attempt 52 the loss was 33.342831, best 30.556467
In attempt 53 the loss was 33.365239, best 30.556467
In attempt 54 the loss was 34.570113, best 30.556467
In attempt 55 the loss was 33.290241, best 30.556467
In attempt 56 the loss was 35.698501, best 30.556467
In attempt 57 the loss was 32.386731, best 30.556467
In attempt 58 the loss was 32.115657, best 30.556467
In attempt 59 the loss was 33.883616, best 30.556467
In attempt 60 the loss was 34.300839, best 30.556467
In attempt 61 the loss was 33.326217, best 30.556467
In attempt 62 the loss was 31.564005, best 30.556467
In attempt 63 the loss was 33.877560, best 30.556467
In attempt 64 the loss was 33.996048, best 30.556467
In attempt 65 the loss was 30.040423, best 30.040423
In attempt 66 the loss was 35.45894, best 30.040423
In attempt 67 the loss was 34.86141, best 30.040423
In attempt 68 the loss was 33.210894, best 30.040423
In attempt 69 the loss was 34.08812, best 30.040423
In attempt 70 the loss was 32.589782, best 30.040423
In attempt 71 the loss was 30.780545, best 30.040423
In attempt 72 the loss was 32.135493, best 30.040423
In attempt 73 the loss was 33.909036, best 30.040423
In attempt 74 the loss was 31.544491, best 30.040423
In attempt 75 the loss was 31.737214, best 30.040423
In attempt 76 the loss was 31.410734, best 30.040423
In attempt 77 the loss was 32.932679, best 30.040423
In attempt 78 the loss was 35.623364, best 30.040423
In attempt 79 the loss was 31.773960, best 30.040423
In attempt 80 the loss was 35.066745, best 30.040423
In attempt 81 the loss was 33.948174, best 30.040423
In attempt 82 the loss was 32.162808, best 30.040423
In attempt 83 the loss was 31.921446, best 30.040423
In attempt 84 the loss was 30.518599, best 30.040423
In attempt 85 the loss was 36.42637, best 30.040423
In attempt 86 the loss was 32.381492, best 30.040423
In attempt 87 the loss was 33.151785, best 30.040423
In attempt 88 the loss was 33.853117, best 30.040423
In attempt 89 the loss was 33.642947, best 30.040423
In attempt 90 the loss was 33.865483, best 30.040423
In attempt 91 the loss was 31.825663, best 30.040423
In attempt 92 the loss was 34.795321, best 30.040423
In attempt 93 the loss was 33.785028, best 30.040423
In attempt 94 the loss was 31.951734, best 30.040423
In attempt 95 the loss was 32.916134, best 30.040423
In attempt 96 the loss was 33.177653, best 30.040423
In attempt 97 the loss was 33.900383, best 30.040423
In attempt 98 the loss was 33.620565, best 30.040423
In attempt 99 the loss was 34.202109, best 30.040423

In [31]: # How bestW perform:
print('Accuracy on train set: ', np.sum(np.argmax(np.abs(1 - X_dev.dot(W)), axis=1) == y_dev).astype(np.float32)/y_train.shape[0])
print('Accuracy on test set: ', np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0])

Accuracy on train set: 10.0
Accuracy on test set: 9.5

You can clearly see that the performance is very low, almost at the random level.
```

### Closed-form solution

The closed-form solution is achieved by:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{W}} &= \frac{2}{n} \mathbf{X}^T (\mathbf{XW} - \mathbf{y}) = 0 \\ \Leftrightarrow \mathbf{W}^* &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

```
In [32]: # TODO:
# Implement the closed-form solution of the weight W.
xtx = np.transpose(X_train) @ X_train
W = np.linalg.inv(xtx).dot(xty)

In [33]: # Check accuracy:
print('Train set accuracy: ', np.sum(np.argmax(np.abs(1 - X_train.dot(W)), axis=1) == y_train).astype(np.float32)/y_train.shape[0])
print('Test set accuracy: ', np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0])

Train set accuracy: 51.163265306122454
Test set accuracy: 36.199999999999996

Now, you can see that the performance is much better.
```

### Regularization

A simple way to improve performance is to include the L2-regularization penalty.

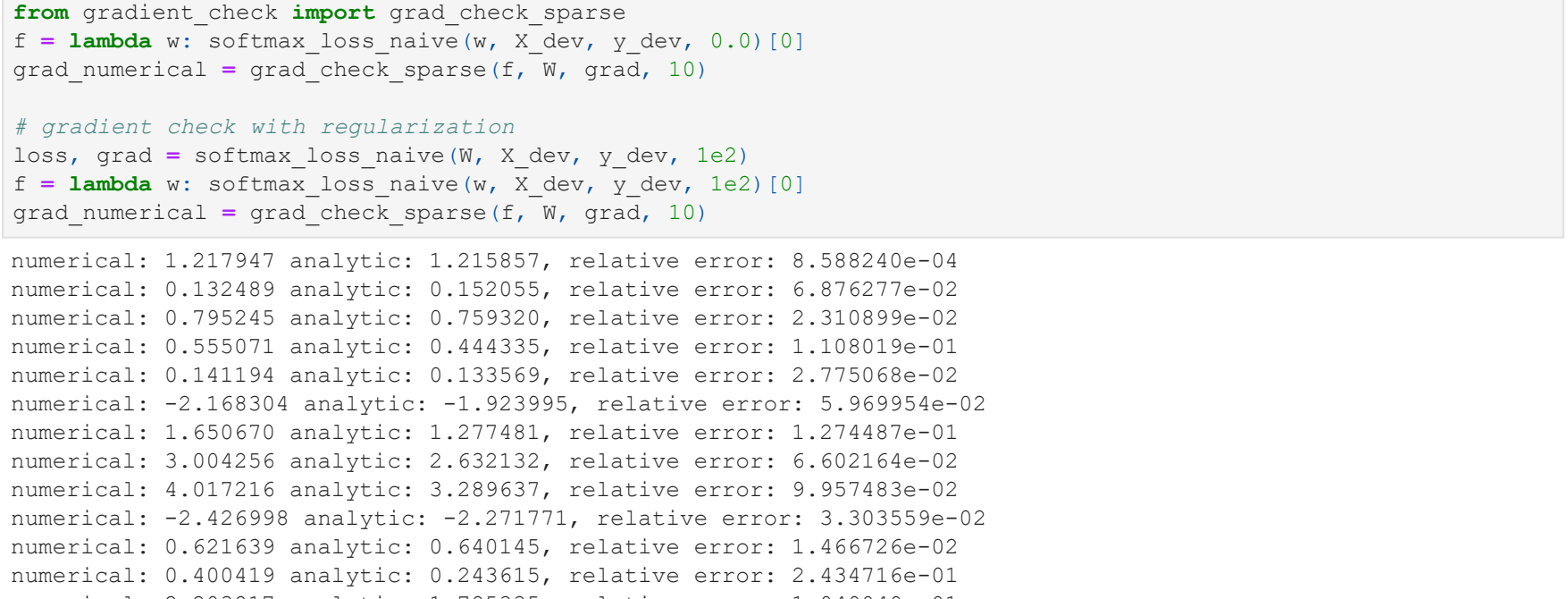
$$\mathcal{L} = \frac{1}{n} \| \mathbf{XW} - \mathbf{y} \|_2^2 + \lambda \| \mathbf{W} \|_2^2 \quad (2)$$

The closed-form solution now is:

$$\Leftrightarrow \mathbf{W}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

```
In [34]: # try several values of lambda to see how it helps:
lambdas = [0.01, 0.1, 1, 10, 100, 1000, 10000, 100000]
train_acc = np.zeros((len(lambdas)))
test_acc = np.zeros((len(lambdas)))
for i in range(len(lambdas)):
    lambda_i = lambdas[i]
    n,d = X_train.shape[0], X_train.shape[1]
    # TODO:
    # Implement the closed-form solution of the weight W with regularization.
    xtx = (np.transpose(X_train) @ X_train) + 1*n*np.eye(d)
    xty = np.transpose(X_train).dot(y_train_oh)
    W = np.linalg.inv(xtx).dot(xty)
    # END OF YOUR CODE

In [35]: plt.semilogx(lambdas, train_acc, 'r', label="Training accuracy")
plt.semilogx(lambdas, test_acc, 'g', label="Testing accuracy")
plt.legend()
plt.grid(True)
plt.show()
```



**Question:** Why to explain why the performances on the training and test set have such behaviors as we change the value of  $\lambda$ .

**Your answer:** As  $\lambda$  increases, a simple model becomes more important than a model that fits the training set very well. Thus, at small values of  $\lambda$ , the training accuracy will be very high and the testing accuracy will be very low as the model will closely fit the training data but not the test data. As  $\lambda$  increases, training accuracy decreases as the model will not fit the training set as well before. Testing accuracy will increase as the model will be more general, but after a certain point,  $\lambda = 1000$  in the above graph, testing accuracy drops as the model becomes too general.

### Softmax Classifier

The predicted probability for the  $j$ -th class given a sample vector  $\mathbf{x}$  and a weight  $\mathbf{W}$  is:

$$P(y = j | \mathbf{x}) = \frac{e^{-\mathbf{xW}_j}}{\sum_{c=1}^C e^{-\mathbf{xW}_c}}$$



Your code for this section will all be written inside `classifiers/softmax.py` and `softmax.py`.

```
In [36]: # First implement the naive softmax loss function with nested loops.
# Open the file classifiers/softmax.py to implement the
# softmax_loss_naive function.

from 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 softmax loss check, our loss should be something close to -log(0.1).
print('loss: %f' % loss)
print('sanity check: %f' % (-np.log(0.1)))

loss: 2.380864
sanity check: 2.302585

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

Your answer: This is as the dataset only has 10 classes, thus the probability of getting the right label for each test example is 0.1. Hence, we expect the loss to be close to -log(0.1).

Optimization

Random search

In [37]: best_loss = float('inf')
for num in range(1001):
    W = np.random.randn(3073, 10) * 0.0001
    loss, _ = softmax_loss_naive(W, X_dev, y_dev, 0.0)
    if loss < best_loss:
        best_loss = loss
        bestW = W
    print('In attempt %d the loss was %f, best %f' % (num, loss, best_loss))

In attempt 0 the loss was 2.346017, best 2.346017
In attempt 1 the loss was 2.323649, best 2.323649
In attempt 2 the loss was 2.346411, best 2.323649
In attempt 3 the loss was 2.357735, best 2.323649
In attempt 4 the loss was 2.31117, best 2.323649
In attempt 5 the loss was 2.323231, best 2.323231
In attempt 6 the loss was 2.330263, best 2.323231
In attempt 7 the loss was 2.363373, best 2.323231
In attempt 8 the loss was 2.252869, best 2.323231
In attempt 9 the loss was 2.438606, best 2.323231
In attempt 10 the loss was 2.360946, best 2.323231
In attempt 11 the loss was 2.437149, best 2.323231
In attempt 12 the loss was 2.393287, best 2.323231
In attempt 13 the loss was 2.343271, best 2.323231
In attempt 14 the loss was 2.39024, best 2.323231
In attempt 15 the loss was 2.425717, best 2.323231
In attempt 16 the loss was 2.392865, best 2.323231
In attempt 17 the loss was 2.357944, best 2.323231
In attempt 18 the loss was 2.37994, best 2.323231
In attempt 19 the loss was 2.357643, best 2.323231
In attempt 20 the loss was 2.316299, best 2.323231
In attempt 21 the loss was 2.334672, best 2.316299
In attempt 22 the loss was 2.404035, best 2.316299
In attempt 23 the loss was 2.304163, best 2.304163
In attempt 24 the loss was 2.386964, best 2.304163
In attempt 25 the loss was 2.339042, best 2.304163
In attempt 26 the loss was 2.41590, best 2.304163
In attempt 27 the loss was 2.427842, best 2.304163
In attempt 28 the loss was 2.327053, best 2.304163
In attempt 29 the loss was 2.375629, best 2.304163
In attempt 30 the loss was 2.346430, best 2.304163
In attempt 31 the loss was 2.314503, best 2.304163
In attempt 32 the loss was 2.372458, best 2.304163
In attempt 33 the loss was 2.380628, best 2.304163
In attempt 34 the loss was 2.360905, best 2.304163
In attempt 35 the loss was 2.307161, best 2.304163
In attempt 36 the loss was 2.399699, best 2.304163
In attempt 37 the loss was 2.349561, best 2.304163
In attempt 38 the loss was 2.435106, best 2.304163
In attempt 39 the loss was 2.382840, best 2.304163
In attempt 40 the loss was 2.311527, best 2.304163
In attempt 41 the loss was 2.297903, best 2.297903
In attempt 42 the loss was 2.315568, best 2.297903
In attempt 43 the loss was 2.370139, best 2.297903
In attempt 44 the loss was 2.321405, best 2.297903
In attempt 45 the loss was 2.350125, best 2.297903
In attempt 46 the loss was 2.314503, best 2.297903
In attempt 47 the loss was 2.388527, best 2.297903
In attempt 48 the loss was 2.323985, best 2.297903
In attempt 49 the loss was 2.37174, best 2.297903
In attempt 50 the loss was 2.352724, best 2.297903
In attempt 51 the loss was 2.353139, best 2.297903
In attempt 52 the loss was 2.471915, best 2.297903
In attempt 53 the loss was 2.435294, best 2.297903
In attempt 54 the loss was 2.368835, best 2.297903
In attempt 55 the loss was 2.343611, best 2.297903
In attempt 56 the loss was 2.370785, best 2.297903
In attempt 57 the loss was 2.413524, best 2.297903
In attempt 58 the loss was 2.340411, best 2.297903
In attempt 59 the loss was 2.306470, best 2.297903
In attempt 60 the loss was 2.324611, best 2.297903
In attempt 61 the loss was 2.278780, best 2.297903
In attempt 62 the loss was 2.439120, best 2.278780
In attempt 63 the loss was 2.361272, best 2.278780
In attempt 64 the loss was 2.383639, best 2.278780
In attempt 65 the loss was 2.294905, best 2.269405
In attempt 66 the loss was 2.411336, best 2.269405
In attempt 67 the loss was 2.396853, best 2.269405
In attempt 68 the loss was 2.380937, best 2.269405
In attempt 69 the loss was 2.394285, best 2.269405
In attempt 70 the loss was 2.302930, best 2.269405
In attempt 71 the loss was 2.356040, best 2.269405
In attempt 72 the loss was 2.373648, best 2.269405
In attempt 73 the loss was 2.370950, best 2.269405
In attempt 74 the loss was 2.309778, best 2.269405
In attempt 75 the loss was 2.341909, best 2.269405
In attempt 76 the loss was 2.311944, best 2.269405
In attempt 77 the loss was 2.362005, best 2.269405
In attempt 78 the loss was 2.406865, best 2.269405
In attempt 79 the loss was 2.369891, best 2.269405
In attempt 80 the loss was 2.339579, best 2.269405
In attempt 81 the loss was 2.294596, best 2.269405
In attempt 82 the loss was 2.353653, best 2.269405
In attempt 83 the loss was 2.371741, best 2.269405
In attempt 84 the loss was 2.335790, best 2.269405
In attempt 85 the loss was 2.323281, best 2.269405
In attempt 86 the loss was 2.372727, best 2.269405
In attempt 87 the loss was 2.345165, best 2.269405
In attempt 88 the loss was 2.340496, best 2.269405
In attempt 89 the loss was 2.413736, best 2.269405
In attempt 90 the loss was 2.276195, best 2.269405
In attempt 91 the loss was 2.335273, best 2.269405
In attempt 92 the loss was 2.364019, best 2.269405
In attempt 93 the loss was 2.336443, best 2.269405
In attempt 94 the loss was 2.363609, best 2.269405
In attempt 95 the loss was 2.343149, best 2.269405
In attempt 96 the loss was 2.323139, best 2.269405
In attempt 97 the loss was 2.333862, best 2.269405
In attempt 98 the loss was 2.401409, best 2.269405
In attempt 99 the loss was 2.375979, best 2.269405

In [38]: # How bestW perform on trainset
scores = X_train.dot(bestW)
y_pred = np.argmax(scores, axis=1)
print('Accuracy on train set: %f' % np.mean(y_pred == y_train))

# evaluate performance of test set
scores = X_test.dot(bestW)
y_pred = np.argmax(scores, axis=1)
print('Accuracy on test set: %f' % np.mean(y_pred == y_test))

Accuracy on train set: 0.136449
Accuracy on test set: 0.143000

Compare the performance when using random search with regression classifier and softmax classifier. You can see how much useful the softmax classifier is.
```

### Stochastic Gradient Descent

Even though it is possible to achieve closed-form solution with softmax classifier, it would be more complicated. In fact, we could achieve very good results with gradient descent approach. Additionally, in case of very large dataset, it is impossible to load the whole dataset into the memory. Gradient descent can help to optimize the loss function in batch.

$$\mathbf{W}^{t+1} = \mathbf{W}^t - \alpha \frac{\partial \mathcal{L}(\mathbf{x}; \mathbf{W}^t)}{\partial \mathbf{W}^t}$$

Where  $\alpha$  is the learning rate,  $\mathcal{L}$  is a loss function, and  $\mathbf{x}$  is a batch of training dataset.

```
In [39]: # 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)

# Use numeric gradient checking as a debugging tool.
# The numeric gradient should be close to the analytic gradient.
from grad_check import grad_check_sparse
f = lambda w: softmax_loss_naive(w, X_dev, y_dev, 0.0)[0]
numerical = grad_check_sparse(f, W, grad, 10)

# gradient check with regularization
loss, grad = softmax_loss_naive(W, X_dev, y_dev, 1e2)
numerical = softmax_loss_naive(W, X_dev, y_dev, 1e2)[0]
numerical = grad_check_sparse(f, W, grad, 10)

numerical: 1.273947 analytic: 1.215857, relative error: 8.568240e-02
numerical: 0.132489 analytic: 0.152055, relative error: 6.876277e-02
numerical: 0.795245 analytic: 0.759320, relative error: 2.310899e-02
numerical: 0.555071 analytic: 0.444335, relative error: 1.108019e-01
numerical: 0.141594 analytic: 0.153663, relative error: 2.775068e-02
numerical: -1.68304 analytic: -1.923995, relative error: 3.969954e-02
numerical: 1.650670 analytic: 1.277481, relative error: 1.274487e-01
numerical: 3.044256 analytic: 2.632132, relative error: 6.602164e-02
numerical: 4.017216 analytic: 3.288637, relative error: 9.957483e-02
numerical: -2.426998 analytic: -2.271771, relative error: 3.303559e-02
numerical: 0.621639 analytic: 0.640145, relative error: 1.466726e-02
numerical: 0.400419 analytic: 0.243615, relative error: 2.474487e-01
numerical: 2.201077 analytic: 2.176535, relative error: 1.049049e-01
numerical: 0.798253 analytic: 0.540171, relative error: 1.920719e-01
numerical: 2.089670 analytic: 1.897814, relative error: 4.811454e-02
numerical: 0.651894 analytic: 0.576449, relative error: 6.602164e-02
numerical: -0.552466 analytic: -0.586683, relative error: 3.003939e-02
numerical: -1.851292 analytic: -1.757381, relative error: 2.602369e-02
numerical: 0.402573 analytic: 0.434034, relative error: 3.760599e-02
numerical: -3.189398 analytic: -2.942604, relative error: 4.024692e-02

In [40]: # 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.
t0 = time.time()
loss_naive, grad_naive = softmax_loss_naive(W, X_dev, y_dev, 0.00001)
t1 = time.time()
print('naive loss: %e computed in %fs' % (loss_naive, t1 - t0))

from classifiers.softmax import softmax_loss_vectorized
t0 = time.time()
loss_vectorized, grad_vectorized = softmax_loss_vectorized(W, X_dev, y_dev, 0.00001)
t1 = time.time()
print('vectorized loss: %e computed in %fs' % (loss_vectorized, t1 - t0))

# 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('Grad difference: %f' % grad_difference)

naive loss: 2.375979e+00 computed in 0.266202s
vectorized loss: 2.375979e+00 computed in 0.001703s
Loss difference: 0.000000
```