

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 validation set to tune the learning rate and regularization strength
- optimize the loss function with SGD
- visualize the final learned weights

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

```
In [1]: import random
import numpy as np
from data_utils import load_CIFAR10
import matplotlib.pyplot as plt

#matplotlib inline
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/1507993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

In [2]: def rel_error(out, correct_out):
    return np.sum(abs(out - correct_out) / (abs(out) + abs(correct_out)))

In [3]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000, num_dev=500):
    """Get the CIFAR-10 data 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 = 'data/cifar10_batch.npy'
    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 will also make a development set, which is a small subset of
    # the training set.
    num_val = 1000
    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,)
```

In [4]: # Create one-hot vectors for labels
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
y_val_oh = np.zeros((y_val.shape[0], 10))
y_val_oh[np.arange(y_val.shape[0]), y_val] = 1
y_test_oh = np.zeros((y_test.shape[0], 10))
y_test_oh[np.arange(y_test.shape[0]), y_test] = 1
y_dev_oh = np.zeros((y_dev.shape[0], 10))
print('Validation accuracy: %f' % (y_dev == 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:

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

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 [5]: bestloss = float('inf')
for num in range(100):
    W = np.random.randn(3073, 10) * 0.0001
    loss = np.linalg.norm(X_dev.dot(W) - y_dev_oh)
    if loss < bestloss:
        bestloss = loss
        bestW = W
    print('in attempt %d the loss was %f, best %f' % (num, loss, bestloss))

in attempt 0 the loss was 31.984276, best 31.984276
in attempt 1 the loss was 33.160889, best 31.984276
in attempt 2 the loss was 33.268427, best 31.984276
in attempt 3 the loss was 31.424596, best 31.424596
in attempt 4 the loss was 34.103124, best 31.424596
in attempt 5 the loss was 32.118090, best 31.424596
in attempt 6 the loss was 37.297972, best 31.424596
in attempt 7 the loss was 35.418770, best 31.424596
in attempt 8 the loss was 32.165426, best 31.424596
in attempt 9 the loss was 36.664584, best 31.424596
in attempt 10 the loss was 31.229600, best 31.229600
in attempt 11 the loss was 31.363957, best 31.229600
in attempt 12 the loss was 33.149059, best 31.229600
in attempt 13 the loss was 32.377586, best 31.229600
in attempt 14 the loss was 30.942296, best 30.942296
in attempt 15 the loss was 31.505083, best 30.942296
in attempt 16 the loss was 32.603060, best 30.942296
in attempt 17 the loss was 32.958431, best 30.942296
in attempt 18 the loss was 32.527920, best 30.942296
in attempt 19 the loss was 32.992204, best 30.942296
in attempt 20 the loss was 36.083025, best 30.942296
in attempt 21 the loss was 32.336069, best 30.942296
in attempt 22 the loss was 34.125836, best 30.942296
in attempt 23 the loss was 32.677125, best 30.942296
in attempt 24 the loss was 32.758899, best 30.942296
in attempt 25 the loss was 32.618524, best 30.942296
in attempt 26 the loss was 34.637209, best 30.942296
in attempt 27 the loss was 32.376546, best 30.942296
in attempt 28 the loss was 33.393951, best 30.942296
in attempt 29 the loss was 31.507072, best 30.942296
in attempt 30 the loss was 32.769004, best 30.942296
in attempt 31 the loss was 33.088100, best 30.942296
in attempt 32 the loss was 32.893127, best 30.942296
in attempt 33 the loss was 32.95452, best 30.942296
in attempt 34 the loss was 32.525147, best 30.942296
in attempt 35 the loss was 34.256572, best 30.942296
in attempt 36 the loss was 34.249349, best 30.942296
in attempt 37 the loss was 32.516597, best 30.942296
in attempt 38 the loss was 31.366712, best 30.942296
in attempt 39 the loss was 32.314209, best 30.942296
in attempt 40 the loss was 35.629302, best 30.942296
in attempt 41 the loss was 32.340249, best 30.942296
in attempt 42 the loss was 34.752109, best 30.942296
in attempt 43 the loss was 33.613475, best 30.942296
in attempt 44 the loss was 34.904040, best 30.942296
in attempt 45 the loss was 31.613980, best 30.942296
in attempt 46 the loss was 32.235462, best 30.942296
in attempt 47 the loss was 34.467053, best 30.942296
in attempt 48 the loss was 32.340224, best 30.942296
in attempt 49 the loss was 32.426964, best 30.942296
in attempt 50 the loss was 32.264906, best 30.942296
in attempt 51 the loss was 32.953152, best 30.942296
in attempt 52 the loss was 32.585514, best 30.942296
in attempt 53 the loss was 35.005060, best 30.942296
in attempt 54 the loss was 35.761892, best 30.942296
in attempt 55 the loss was 32.265779, best 30.942296
in attempt 56 the loss was 34.692084, best 30.942296
in attempt 57 the loss was 32.264906, best 30.942296
in attempt 58 the loss was 33.953152, best 30.942296
in attempt 59 the loss was 32.608348, best 30.942296
in attempt 60 the loss was 32.204505, best 30.942296
in attempt 61 the loss was 32.764100, best 30.942296
in attempt 62 the loss was 32.666877, best 30.942296
in attempt 63 the loss was 36.719779, best 30.942296
in attempt 64 the loss was 34.762724, best 30.942296
in attempt 65 the loss was 32.933152, best 30.942296
in attempt 66 the loss was 31.921522, best 30.942296
in attempt 67 the loss was 32.933152, best 30.942296
in attempt 68 the loss was 32.518763, best 30.942296
in attempt 69 the loss was 32.305250, best 30.942296
in attempt 70 the loss was 34.421365, best 30.942296
in attempt 71 the loss was 32.393463, best 30.942296
in attempt 72 the loss was 32.75448, best 30.942296
in attempt 73 the loss was 34.338005, best 30.942296
in attempt 74 the loss was 32.933152, best 30.942296
in attempt 75 the loss was 31.549278, best 30.942296
in attempt 76 the loss was 32.60372, best 30.942296
in attempt 77 the loss was 32.229605, best 30.942296
in attempt 78 the loss was 32.740089, best 30.942296
in attempt 79 the loss was 34.064465, best 30.942296
in attempt 80 the loss was 31.873032, best 30.942296
in attempt 81 the loss was 33.493229, best 30.942296
in attempt 82 the loss was 32.244177, best 30.942296
in attempt 83 the loss was 32.907723, best 30.942296
in attempt 84 the loss was 32.544296, best 30.942296
in attempt 85 the loss was 32.342074, best 30.942296
in attempt 86 the loss was 32.474741, best 30.942296
in attempt 87 the loss was 32.369335, best 30.942296
in attempt 88 the loss was 31.600303, best 30.942296
in attempt 89 the loss was 32.549759, best 30.942296
in attempt 90 the loss was 32.259604, best 30.942296
in attempt 91 the loss was 32.953279, best 30.942296
in attempt 92 the loss was 32.719566, best 30.942296
in attempt 93 the loss was 32.495519, best 30.942296
in attempt 94 the loss was 32.688741, best 30.942296
in attempt 95 the loss was 35.348429, best 30.942296
in attempt 96 the loss was 32.618786, best 30.942296
in attempt 97 the loss was 34.527502, best 30.942296
in attempt 98 the loss was 31.373089, best 30.942296
in attempt 99 the loss was 33.74492, best 30.942296

In [6]: # Now bestW perform:
train_acc = np.sum(np.argmax(np.abs(1 - X_train.dot(W)), axis=1) == y_train).astype(np.float32)/y_train.shape[0]*100
print('Accuracy on train set: ', train_acc)
test_acc = np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0]*100
print('Accuracy on test set: ', test_acc)
Accuracy on train set: 9.6
Accuracy on test set: 8.5

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:

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

```
In [7]: # TODO:
# Implement the closed-form solution of the weight W.
#
# no matrix multiplication
A_inv = np.linalg.pinv(A)
B = np.matmul(np.transpose(X_train), y_train_oh)
W = np.matmul(A_inv, B)

In [8]: # 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]*100)
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]*100)
Train set accuracy: 51.63565366122454
Test set accuracy: 31.109609090909090

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

Regularization

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

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

The closed-form solution now is:

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

```
In [9]: # 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, lambda in enumerate(lambdas):
    # no matrix multiplication
    A_inv = np.linalg.pinv(A)
    B = np.matmul(np.transpose(X_train), y_train_oh)
    W = np.matmul(A_inv, B)

    # Implement the closed-form solution of the weight W with regularization.
    #
    # no matrix multiplication
    A_reg = np.matmul(np.transpose(X_train), X_train) + 1 * lambda * np.identity(d)
    A_inv_reg = np.linalg.pinv(A_reg)
    B = np.matmul(np.transpose(X_train), y_train_oh)
    W = np.matmul(A_inv_reg, B)

    # Check accuracy:
    train_acc[i] = np.sum(np.argmax(np.abs(1 - X_train.dot(W)), axis=1) == y_train).astype(np.float32)/y_train.shape[0]*100
    test_acc[i] = np.sum(np.argmax(np.abs(1 - X_test.dot(W)), axis=1) == y_test).astype(np.float32)/y_test.shape[0]*100

In [10]: 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: Try to explain why the performances on the training and test set have such behaviors as we change the value of λ .

Your answer: This is because as the value of lambda gets bigger, the model becomes simpler to prevent overfitting which can be observed as the performance on the test set improves and the performance on training set drops. However, when the lambda value is too high, the model starts to underfit the data resulting in poor performance on both the training and test sets.

Softmax Classifier

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

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

```
%softmax

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

In [11]: # First implement the naive softmax loss function with nested loops.
# Open the file classifiers/softmax.py and 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 sanity check, our loss should be something close to -log(0.1).
print('sanity check: %f' % (-np.log(0.1)))

loss: 2.401973
sanity check: 2.302585

Question: Why do we expect our loss to be close to -log(0.1)? Explain briefly.
Your answer: There are 10 classes so on average there is a 0.1 probability of predicting the correct class. Therefore, we expect loss to be close to -log(0.1)

Optimization

Random search

In [12]: bestloss = float('inf')
for num in range(100):
    W = np.random.randn(3073, 10) * 0.0001
    loss, _ = softmax_loss_naive(W, X_dev, y_dev, 0.0)
    if loss < bestloss:
        bestloss = loss
        bestW = W
    print('in attempt %d the loss was %f, best %f' % (num, loss, bestloss))

in attempt 0 the loss was 2.353272, best 2.353272
in attempt 1 the loss was 2.307397, best 2.307397
in attempt 2 the loss was 2.347314, best 2.307397
in attempt 3 the loss was 2.387799, best 2.307397
in attempt 4 the loss was 2.326592, best 2.307397
in attempt 5 the loss was 2.352080, best 2.307397
in attempt 6 the loss was 2.382080, best 2.307397
in attempt 7 the loss was 2.308069, best 2.307397
in attempt 8 the loss was 2.421393, best 2.307397
in attempt 9 the loss was 2.328353, best 2.307397
in attempt 10 the loss was 2.363496, best 2.307397
in attempt 11 the loss was 2.384060, best 2.307397
in attempt 12 the loss was 2.324681, best 2.307397
in attempt 13 the loss was 2.313397, best 2.307397
in attempt 14 the loss was 2.329612, best 2.307397
in attempt 15 the loss was 2.348474, best 2.307397
in attempt 16 the loss was 2.305164, best 2.307397
in attempt 17 the loss was 2.358240, best 2.305164
in attempt 18 the loss was 2.327029, best 2.305164
in attempt 19 the loss was 2.402366, best 2.305164
in attempt 20 the loss was 2.303051, best 2.305164
in attempt 21 the loss was 2.402366, best 2.305164
in attempt 22 the loss was 2.309352, best 2.305164
in attempt 23 the loss was 2.353957, best 2.305164
in attempt 24 the loss was 2.329169, best 2.305164
in attempt 25 the loss was 2.306254, best 2.305164
in attempt 26 the loss was 2.326852, best 2.305164
in attempt 27 the loss was 2.367389, best 2.305164
in attempt 28 the loss was 2.324909, best 2.305164
in attempt 29 the loss was 2.369748, best 2.305164
in attempt 30 the loss was 2.345002, best 2.305164
in attempt 31 the loss was 2.351576, best 2.305164
in attempt 32 the loss was 2.316730, best 2.305164
in attempt 33 the loss was 2.349269, best 2.305164
in attempt 34 the loss was 2.365370, best 2.305164
in attempt 35 the loss was 2.349185, best 2.305164
in attempt 36 the loss was 2.401383, best 2.305164
in attempt 37 the loss was 2.327120, best 2.305164
in attempt 38 the loss was 2.373833, best 2.305164
in attempt 39 the loss was 2.336420, best 2.305164
in attempt 40 the loss was 2.364465, best 2.305164
in attempt 41 the loss was 2.334910, best 2.305164
in attempt 42 the loss was 2.369889, best 2.305164
in attempt 43 the loss was 2.463806, best 2.305164
in attempt 44 the loss was 2.309495, best 2.305164
in attempt 45 the loss was 2.310923, best 2.305164
in attempt 46 the loss was 2.413060, best 2.305164
in attempt 47 the loss was 2.329550, best 2.305164
in attempt 48 the loss was 2.338083, best 2.305164
in attempt 49 the loss was 2.324331, best 2.305164
in attempt 50 the loss was 2.350986, best 2.305164
in attempt 51 the loss was 2.321705, best 2.305164
in attempt 52 the loss was 2.329749, best 2.305164
in attempt 53 the loss was 2.317668, best 2.305164
in attempt 54 the loss was 2.306669, best 2.305164
in attempt 55 the loss was 2.370277, best 2.305164
in attempt 56 the loss was 2.304090, best 2.305164
in attempt 57 the loss was 2.310973, best 2.305164
in attempt 58 the loss was 2.358957, best 2.305164
in attempt 59 the loss was 2.355775, best 2.305164
in attempt 60 the loss was 2.308097, best 2.305164
in attempt 61 the loss was 2.303206, best 2.305164
in attempt 62 the loss was 2.336293, best 2.305164
in attempt 63 the loss was 2.361712, best 2.305164
in attempt 64 the loss was 2.306935, best 2.305164
in attempt 65 the loss was 2.354545, best 2.305164
in attempt 66 the loss was 2.365893, best 2.305164
in attempt 67 the loss was 2.359595, best 2.305164
in attempt 68 the loss was 2.388693, best 2.305164
in attempt 69 the loss was 2.310183, best 2.305164
in attempt 70 the loss was 2.309609, best 2.305164
in attempt 71 the loss was 2.370485, best 2.305164
in attempt 72 the loss was 2.318649, best 2.305164
in attempt 73 the loss was 2.327455, best 2.305164
in attempt 74 the loss was 2.353406, best 2.305164
in attempt 75 the loss was 2.340423, best 2.305164
in attempt 76 the loss was 2.343052, best 2.305164
in attempt 77 the loss was 2.313649, best 2.305164
in attempt 78 the loss was 2.362760, best 2.305164
in attempt 79 the loss was 2.329417, best 2.305164
in attempt 80 the loss was 2.339333, best 2.290417
in attempt 81 the loss was 2.344500, best 2.290417
in attempt 82 the loss was 2.356484, best 2.290417
in attempt 83 the loss was 2.336295, best 2.290417
in attempt 84 the loss was 2.313695, best 2.290417
in attempt 85 the loss was 2.363995, best 2.290417
in attempt 86 the loss was 2.373203, best 2.290417
in attempt 87 the loss was 2.330871, best 2.290417
in attempt 88 the loss was 2.351447, best 2.290417
in attempt 89 the loss was 2.334307, best 2.290417
in attempt 90 the loss was 2.332400, best 2.290417
in attempt 91 the loss was 2.354545, best 2.290417
in attempt 92 the loss was 2.320203, best 2.290417
in attempt 93 the loss was 2.405210, best 2.290417
in attempt 94 the loss was 2.337899, best 2.290417
in attempt 95 the loss was 2.339375, best 2.290417
in attempt 96 the loss was 2.361363, best 2.290417
in attempt 97 the loss was 2.373476, best 2.290417
in attempt 98 the loss was 2.377789, best 2.290417
in attempt 99 the loss was 2.377789, best 2.290417

In [13]: # Now 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: 0.135480
Accuracy on test set: 0.130800

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's 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.

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

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

```
In [14]: # 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 on a debugging weight.
def gradient_check(x, y, W, grad):
    """The numeric gradient should be close to the analytic gradient.
    """
    from gradient_check import grad_check_sparse
    (lambda_w, softmax_loss_naive(W, X_dev, y_dev, 0.0))
    grad_numerical = grad_check_sparse(W, grad, 10)

    # gradient check with regularization
    loss, grad = softmax_loss_naive(W, X_dev, y_dev, 1e2)
    grad_numerical = grad_check_sparse(W, grad, 10)

    numerical: 2.397732 analytic: 34.552780, relative error: 9.607843e-01
    numerical: 0.691056 analytic: 34.552780, relative error: 9.607843e-01
    numerical: 0.519029 analytic: 25.908995, relative error: 9.607843e-01
    numerical: 2.319177 analytic: 119.858849, relative error: 9.607843e-01
    numerical: 0.072941 analytic: 3.647842, relative error: 9.607843e-01
    numerical: 0.951304 analytic: 1.814971, relative error: 9.607843e-01
    numerical: -2.018192 analytic: -104.050113, relative error: 9.607843e-01
    numerical: 0.859718 analytic: 42.959511, relative error: 9.607843e-01
    numerical: 1.582886 analytic: 71.191239, relative error: 9.607843e-01
    numerical: 1.436248 analytic: 71.812382, relative error: 9.607843e-01
    numerical: 1.987069 analytic: 53.927239, relative error: 9.607843e-01
    numerical: -4.229027 analytic: 120.594553, relative error: 9.609131e-01
    numerical: 1.764479 analytic: 85.088837, relative error: 9.607196e-01
    numerical: 2.563051 analytic: 120.224206, relative error: 9.608017e-01
    numerical: -0.170297 analytic: -0.788768, relative error: 9.619875e-01
    numerical: -0.125851 analytic: -0.518551, relative error: 9.556406e-01
    numerical: -0.142088 analytic: -0.648971, relative error: 9.581106e-01
    numerical: -1.225077 analytic: -61.856828, relative error: 9.611502e-01
    numerical: 1.626986 analytic: 61.231304, relative error: 9.607209e-01

In [15]: # 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, 0.00001)
toc = time.time()
print('naive loss: %f computed in %fs' % (loss_naive, toc - tic))

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

# 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' % abs(loss_naive - loss_vectorized))
print('gradient difference: %f' % grad_difference)

naive loss: 2.377789e+00 computed in 0.009755s
vectorized loss: 2.377789e+00 computed in 0.002355s
loss difference: 0.000000
gradient difference: 0.000000

In [16]: from classifiers.linear_classifier import *

classifier = Softmax()
tic = time.time()
y_train_pred = classifier.train(X_train, y_train, learning_rate=1e-7, reg=5e-4, num_iters=1500, verbose=True)
toc = time.time()
print('That took %fs' % (toc - tic))

iteration 
```