## This is the k-nearest neighbors workbook for ECE 239AS Assignment #2

Please follow the notebook linearly to implement k-nearest neighbors.

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and code in the jupyer notebook to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

The goal of this workbook is to give you experience with the data, training and evaluating a simple classifier, k-fold cross validation, and as a Python refresher.

## Import the appropriate libraries

```
In [1]: import numpy as np # for doing most of our calculations
        import matplotlib.pyplot as plt# for plotting
        from cs231n.data_utils import load_CIFAR10 # function to load the CIFAR-10
        dataset.

        # Load matplotlib images inline
        %matplotlib inline

        # These are important for reloading any code you write in external .py file
        s.
        # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-i
        python
        %load_ext autoreload
        %autoreload 2
```

```
In [2]: # Set the path to the CIFAR-10 data
        cifar10_dir = 'cifar-10-batches-py'
        X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

        # As a sanity check, we print out the size of the training and test data.
        print('Training data shape: ', X_train.shape)
        print('Training labels shape: ', y_train.shape)
        print('Test data shape: ', X_test.shape)
        print('Test labels shape: ', y_test.shape)
```

```
('Training data shape: ', (50000, 32, 32, 3))
('Training labels shape: ', (50000,))
('Test data shape: ', (10000, 32, 32, 3))
('Test labels shape: ', (10000,))
```

In [3]:
```python
# Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', '
ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()
```



In [4]:
```python
# Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]

num_test = 500
mask = list(range(num_test))
X_test = X_test[mask]
y_test = y_test[mask]

# Reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
print(X_train.shape, X_test.shape)
```

((5000, 3072), (500, 3072))

## K-nearest neighbors

In the following cells, you will build a KNN classifier and choose hyperparameters via k-fold cross-validation.

```
In [5]:  # Import the KNN class

         from nndl import KNN
```

```
In [6]:  # Declare an instance of the knn class.
         knn = KNN()

         # Train the classifier.
         #   We have implemented the training of the KNN classifier.
         #   Look at the train function in the KNN class to see what this does.
         knn.train(X=X_train, y=y_train)
```

## Questions

(1) Describe what is going on in the function knn.train().

(2) What are the pros and cons of this training step?

## Answers

(1) The knn object is copying all of the training data into internal memory so that the points can be used as neighbors in the classifier.

(2)

```
    * Pros
      * Very simple training algorithm, takes only two lines of code
      * O(1) to copy the memory
    * Cons
      * Expensive memory wise
```

## KNN prediction

In the following sections, you will implement the functions to calculate the distances of test points to training points, and from this information, predict the class of the KNN.

In [7]:
```
# Implement the function compute_distances() in the KNN class.
# Do not worry about the input 'norm' for now; use the default definition o
f the norm
#   in the code, which is the 2-norm.
# You should only have to fill out the clearly marked sections.

import time
time_start =time.time()

dists_L2 = knn.compute_distances(X=X_test)
print('Time to run code: {}'.format(time.time()-time_start))
print('Frobenius norm of L2 distances: {}'.format(np.linalg.norm(dists_L2,
'fro')))
```

```
Time to run code: 30.6747307777
Frobenius norm of L2 distances: 7906696.07704
```

**Really slow code**

Note: This probably took a while. This is because we use two for loops. We could increase the speed via vectorization, removing the for loops.

If you implemented this correctly, evaluating np.linalg.norm(dists_L2, 'fro') should return: ~7906696

## KNN vectorization

The above code took far too long to run. If we wanted to optimize hyperparameters, it would be time-expensive. Thus, we will speed up the code by vectorizing it, removing the for loops.

In [8]:
```
# Implement the function compute_L2_distances_vectorized() in the KNN class
.
# In this function, you ought to achieve the same L2 distance but WITHOUT a
ny for loops.
# Note, this is SPECIFIC for the L2 norm.

time_start =time.time()
dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_test)
print('Time to run code: {}'.format(time.time()-time_start))
print('Difference in L2 distances between your KNN implementations (should
be 0): {}'.format(np.linalg.norm(dists_L2 - dists_L2_vectorized, 'fro')))
```

```
Time to run code: 0.22722697258
Difference in L2 distances between your KNN implementations (should be 0):
0.0
```

**Speedup**

Depending on your computer speed, you should see a 10-100x speed up from vectorization. On our computer, the vectorized form took 0.36 seconds while the naive implementation took 38.3 seconds.

### Implementing the prediction

Now that we have functions to calculate the distances from a test point to given training points, we now implement the function that will predict the test point labels.

```
In [9]: # Implement the function predict_labels in the KNN class.
        # Calculate the training error (num_incorrect / total_samples)
        #    from running knn.predict_labels with k=1

        error = 1
        # ================================================================= #
        # YOUR CODE HERE:
        #    Calculate the error rate by calling predict_labels on the test
        #    data with k = 1.  Store the error rate in the variable error.
        # ================================================================= #

        pLabels = knn.predict_labels(dists_L2_vectorized, k=1)

        errors = pLabels - y_test
        errors[np.nonzero(errors)] = 1

        error = sum(errors)/num_test



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

        print(error)
```

```
0.726
```

If you implemented this correctly, the error should be: 0.726.

This means that the k-nearest neighbors classifier is right 27.4% of the time, which is not great, considering that chance levels are 10%.

# Optimizing KNN hyperparameters

In this section, we'll take the KNN classifier that you have constructed and perform cross-validation to choose a best value of $k$, as well as a best choice of norm.

### Create training and validation folds

First, we will create the training and validation folds for use in k-fold cross validation.

```
In [10]:  # Create the dataset folds for cross-valdiation.
          num_folds = 5

          X_train_folds = []
          y_train_folds =  []

          # ================================================================ #
          # YOUR CODE HERE:
          #   Split the training data into num_folds (i.e., 5) folds.
          #   X_train_folds is a list, where X_train_folds[i] contains the
          #       data points in fold i.
          #   y_train_folds is also a list, where y_train_folds[i] contains
          #       the corresponding labels for the data in X_train_folds[i]
          # ================================================================ #

          fold_size = num_training/num_folds
          print X_train.shape
          print y_train.shape

          perm = np.arange(X_train.shape[0])
          np.random.shuffle(perm)
          X_train_shuffled = X_train[perm,:]
          y_train_shuffled = y_train[perm]

          for i in range(num_folds):

              X_train_folds.append(X_train_shuffled[i*(fold_size):(i+1)*(fold_size)])
              y_train_folds.append(y_train_shuffled[i*(fold_size):(i+1)*(fold_size)])

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

```
(5000, 3072)
(5000,)
```

### Optimizing the number of nearest neighbors hyperparameter.

In this section, we select different numbers of nearest neighbors and assess which one has the lowest k-fold cross validation error.

```python
In [11]: time_start =time.time()

         ks = [1, 2, 3, 5, 7, 10, 15, 20, 25, 30]

         # ================================================================ #
         # YOUR CODE HERE:
         #   Calculate the cross-validation error for each k in ks, testing
         #   the trained model on each of the 5 folds.  Average these errors
         #   together and make a plot of k vs. cross-validation error. Since
         #   we are assuming L2 distance here, please use the vectorized code!
         #   Otherwise, you might be waiting a long time.
         # ================================================================ #

         for i in np.arange(len(ks)):

             error = 0

             for j in np.arange(num_folds):

                 X_validate = X_train_folds[j]
                 y_validate = y_train_folds[j]

                 X_training = np.vstack(X_train_folds[:j] + X_train_folds[j+1:])
                 y_training = np.concatenate(y_train_folds[:j] + y_train_folds[j+1:]
         )

                 knn.train(X=X_training, y=y_training)

                 dists = knn.compute_L2_distances_vectorized(X=X_validate)

                 pLabels = knn.predict_labels(dists, k=ks[i])

                 errors = pLabels - y_validate
                 errors[np.nonzero(errors)] = 1

                 error += sum(errors)/fold_size

             error = error/num_folds
             print error


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

         print('Computation time: %.2f'%(time.time()-time_start))
```

```
0.7315999999999999
0.7602
0.744
0.7285999999999999
0.7262000000000001
0.7302
0.7294
0.7242
0.7258
0.7327999999999999
Computation time: 33.87
```

## Questions:

(1) What value of $k$ is best amongst the tested $k$'s?

(2) What is the cross-validation error for this value of $k$?

## Answers:

(1) A k of 20 showed the least cross-validation error

(2) The error was .7242

### Optimizing the norm

Next, we test three different norms (the 1, 2, and infinity norms) and see which distance metric results in the best cross-validation performance.

In [12]:
```python
time_start =time.time()

L1_norm = lambda x: np.linalg.norm(x, ord=1)
L2_norm = lambda x: np.linalg.norm(x, ord=2)
Linf_norm = lambda x: np.linalg.norm(x, ord= np.inf)
norms = [L1_norm, L2_norm, Linf_norm]

# ================================================================ #
# YOUR CODE HERE:
#    Calculate the cross-validation error for each norm in norms, testing
#    the trained model on each of the 5 folds.  Average these errors
#    together and make a plot of the norm used vs the cross-validation error
#    Use the best cross-validation k from the previous part.
#
#    Feel free to use the compute_distances function.  We're testing just
#    three norms, but be advised that this could still take some time.
#    You're welcome to write a vectorized form of the L1- and Linf- norms
#    to speed this up, but it is not necessary.
# ================================================================ #
for i in np.arange(len(norms)):

    error = 0

    for j in np.arange(num_folds):

        X_validate = X_train_folds[j]
        y_validate = y_train_folds[j]

        X_training = np.vstack(X_train_folds[:j] + X_train_folds[j+1:])
        y_training = np.concatenate(y_train_folds[:j] + y_train_folds[j+1:]
)

        knn.train(X=X_training, y=y_training)

        dists = knn.compute_distances(X=X_validate, norm=norms[i])

        pLabels = knn.predict_labels(dists, k=20)

        errors = pLabels - y_validate
        errors[np.nonzero(errors)] = 1

        error += sum(errors)/fold_size

    error = error/num_folds
    print error

# ================================================================ #
# END YOUR CODE HERE
# ================================================================ #
print('Computation time: %.2f'%(time.time()-time_start))
```

```
0.6988
0.7242
0.8321999999999999
Computation time: 759.01
```

## Questions:

(1) What norm has the best cross-validation error?

(2) What is the cross-validation error for your given norm and k?

## Answers:

(1) The L1 norm has the best cross-validation error

(2) It showed a cross-validation error of .6988

# Evaluating the model on the testing dataset.

Now, given the optimal $k$ and norm you found in earlier parts, evaluate the testing error of the k-nearest neighbors model.

```
In [13]: error = 1

# ================================================================ #
# YOUR CODE HERE:
#   Evaluate the testing error of the k-nearest neighbors classifier
#   for your optimal hyperparameters found by 5-fold cross-validation.
# ================================================================ #

knn.train(X=X_train, y=y_train)

dists = knn.compute_distances(X=X_test,norm=L1_norm)

pLabels = knn.predict_labels(dists, k=20)

errors = pLabels - y_test
errors[np.nonzero(errors)] = 1

error = sum(errors)/num_test




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

print('Error rate achieved: {}'.format(error))

Error rate achieved: 0.72
```

## Question:

How much did your error improve by cross-validation over naively choosing $k = 1$ and using the L2-norm?

## Answer:

Only around .6%

```python
1   import numpy as np
2   import pdb
3
4   """
5   This code was based off of code from cs231n at Stanford University, and modified for ece239as at UCLA.
6   """
7
8   class KNN(object):
9
10    def __init__(self):
11      pass
12
13    def train(self, X, y):
14      """
15    Inputs:
16    - X is a numpy array of size (num_examples, D)
17    - y is a numpy array of size (num_examples, )
18      """
19      self.X_train = X
20      self.y_train = y
21
22    def compute_distances(self, X, norm=None):
23      """
24      Compute the distance between each test point in X and each training point
25      in self.X_train.
26
27      Inputs:
28      - X: A numpy array of shape (num_test, D) containing test data.
29    - norm: the function with which the norm is taken.
30
31      Returns:
32      - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
33        is the Euclidean distance between the ith test point and the jth training
34        point.
35      """
36      if norm is None:
37        norm = lambda x: np.sqrt(np.sum(x**2))
38        #norm = 2
39
40      num_test = X.shape[0]
41      num_train = self.X_train.shape[0]
42      dists = np.zeros((num_test, num_train))
43      for i in np.arange(num_test):
44
45        for j in np.arange(num_train):
46          # ================================================================ #
47          # YOUR CODE HERE:
48          #   Compute the distance between the ith test point and the jth
49          #   training point using norm(), and store the result in dists[i, j].
50          # ================================================================ #
51
52          dists[i,j] = norm(X[i] - self.X_train[j])
53
54          # ================================================================ #
55          # END YOUR CODE HERE
56          # ================================================================ #
57
58      return dists
59
60    def compute_L2_distances_vectorized(self, X):
61      """
62      Compute the distance between each test point in X and each training point
63      in self.X_train WITHOUT using any for loops.
64
65      Inputs:
66      - X: A numpy array of shape (num_test, D) containing test data.
67
```

```python
68        Returns:
69        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
70          is the Euclidean distance between the ith test point and the jth training
71          point.
72        """
73        num_test = X.shape[0]
74        #print 'x shape ' , X.shape
75        num_train = self.X_train.shape[0]
76        #print 'x train shape ' , self.X_train.shape
77        dists = np.zeros((num_test, num_train))
78
79        # ================================================================ #
80        # YOUR CODE HERE:
81        #   Compute the L2 distance between the ith test point and the jth
82        #   training point and store the result in dists[i, j].  You may
83        #    NOT use a for loop (or list comprehension).  You may only use
84        #     numpy operations.
85        #
86        #     HINT: use broadcasting.  If you have a shape (N,1) array and
87        #   a shape (M,) array, adding them together produces a shape (N, M)
88        #    array.
89        # ================================================================ #
90
91        xsquared = np.sum(X**2, axis=1)[:, np.newaxis]
92        x_trainsquared = np.sum(self.X_train**2, axis=1)
93        xdotx_train = -2*np.dot(X,self.X_train.T)
94
95
96        dists = np.sqrt(xsquared + x_trainsquared + xdotx_train)
97
98        # ================================================================ #
99        # END YOUR CODE HERE
100        # ================================================================ #
101
102        return dists
103
104
105    def predict_labels(self, dists, k=1):
106        """
107        Given a matrix of distances between test points and training points,
108        predict a label for each test point.
109
110        Inputs:
111        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
112          gives the distance betwen the ith test point and the jth training point.
113
114        Returns:
115        - y: A numpy array of shape (num_test,) containing predicted labels for the
116          test data, where y[i] is the predicted label for the test point X[i].
117        """
118        num_test = dists.shape[0]
119        y_pred = np.zeros(num_test)
120        for i in np.arange(num_test):
121          # A list of length k storing the labels of the k nearest neighbors to
122          # the ith test point.
123          closest_y = []
124          # ================================================================ #
125          # YOUR CODE HERE:
126          #   Use the distances to calculate and then store the labels of
127          #   the k-nearest neighbors to the ith test point.  The function
128          #   numpy.argsort may be useful.
129          #
130          #   After doing this, find the most common label of the k-nearest
131          #   neighbors.  Store the predicted label of the ith training example
132          #   as y_pred[i].  Break ties by choosing the smaller label.
133          # ================================================================ #
134
135
```

```
136        sortedIdxs = np.argsort(dists[i,:])[:k]
137        closest_y = self.y_train[sortedIdxs]
138        counts = np.bincount(closest_y)
139        y_pred[i] = np.argmax(counts)
140
141
142        # ================================================================ #
143        # END YOUR CODE HERE
144        # ================================================================ #
145
146    return y_pred
```

## This is the svm workbook for ECE 239AS Assignment #2

Please follow the notebook linearly to implement a linear support vector machine.

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and includes code to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

The goal of this workbook is to give you experience with training an SVM classifier via gradient descent.

## Importing libraries and data setup

```
In [1]:  import numpy as np # for doing most of our calculations
         import matplotlib.pyplot as plt# for plotting
         from cs231n.data_utils import load_CIFAR10 # function to load the CIFAR-10
         dataset.
         import pdb

         # Load matplotlib images inline
         %matplotlib inline

         # These are important for reloading any code you write in external .py file
         s.
         # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-i
         python
         %load_ext autoreload
         %autoreload 2
```

```
In [2]:  # Set the path to the CIFAR-10 data
         cifar10_dir = 'cifar-10-batches-py'
         X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

         # As a sanity check, we print out the size of the training and test data.
         print('Training data shape: ', X_train.shape)
         print('Training labels shape: ', y_train.shape)
         print('Test data shape: ', X_test.shape)
         print('Test labels shape: ', y_test.shape)
```

```
('Training data shape: ', (50000, 32, 32, 3))
('Training labels shape: ', (50000,))
('Test data shape: ', (10000, 32, 32, 3))
('Test labels shape: ', (10000,))
```

In [3]:
```python
# Visualize some examples from the dataset.
# We show a few examples of training images from each class.
classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', '
ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()
```

In [4]:
```python
# Split the data into train, val, and test sets. In addition we will
# create a small development set as a subset of the training data;
# we can use this for development so our code runs faster.
num_training = 49000
num_validation = 1000
num_test = 1000
num_dev = 500

# Our validation set will be num_validation points from the original
# training set.
mask = range(num_training, num_training + num_validation)
X_val = X_train[mask]
y_val = y_train[mask]

# Our training set will be the first num_train points from the original
# training set.
mask = range(num_training)
X_train = X_train[mask]
y_train = y_train[mask]

# We will also make a development set, which is a small subset of
# the training set.
mask = np.random.choice(num_training, num_dev, replace=False)
X_dev = X_train[mask]
y_dev = y_train[mask]

# We use the first num_test points of the original test set as our
# test set.
mask = range(num_test)
X_test = X_test[mask]
y_test = y_test[mask]

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, 32, 32, 3))
('Train labels shape: ', (49000,))
('Validation data shape: ', (1000, 32, 32, 3))
('Validation labels shape: ', (1000,))
('Test data shape: ', (1000, 32, 32, 3))
('Test labels shape: ', (1000,))
('Dev data shape: ', (500, 32, 32, 3))
('Dev labels shape: ', (500,))
```
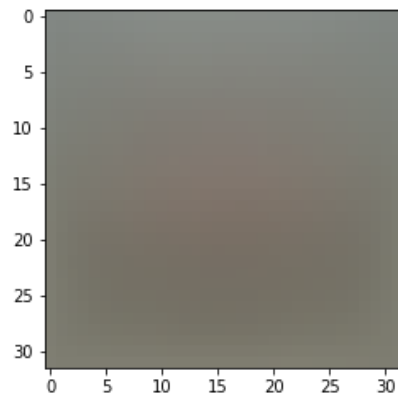
In [5]:
```python
# Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

# As a sanity check, print out the shapes of the data
print('Training data shape: ', X_train.shape)
print('Validation data shape: ', X_val.shape)
print('Test data shape: ', X_test.shape)
print('dev data shape: ', X_dev.shape)
```

```
('Training data shape: ', (49000, 3072))
('Validation data shape: ', (1000, 3072))
('Test data shape: ', (1000, 3072))
('dev data shape: ', (500, 3072))
```

In [6]:
```python
# Preprocessing: subtract the mean image
# first: compute the image mean based on the training data
mean_image = np.mean(X_train, axis=0)
print(mean_image[:10]) # print a few of the elements
plt.figure(figsize=(4,4))
plt.imshow(mean_image.reshape((32,32,3)).astype('uint8')) # visualize the m
ean image
plt.show()
```

```
[130.64189796 135.98173469 132.47391837 130.05569388 135.34804082
 131.75402041 130.96055102 136.14328571 132.47636735 131.48467347]
```



In [7]:
```python
# second: subtract the mean image from train and test data
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image
```

In [8]:
```
# third: append the bias dimension of ones (i.e. bias trick) so that our SV
M
# only has to worry about optimizing a single weight matrix W.
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))])

print(X_train.shape, X_val.shape, X_test.shape, X_dev.shape)
```

((49000, 3073), (1000, 3073), (1000, 3073), (500, 3073))

## Question:

(1) For the SVM, we perform mean-subtraction on the data. However, for the KNN notebook, we did not. Why?

## Answer:

(1) Mean subtracting on KNN would shift all of the points an equal amount, changing nothing about the neighbors. In the SVM we need to mean subtract so that the data is zero centered and the gradients behave well in multiple directions.

## Training an SVM

The following cells will take you through building an SVM. You will implement its loss function, then subsequently train it with gradient descent. Finally, you will choose the learning rate of gradient descent to optimize its classification performance.

In [9]:
```
from nndl.svm import SVM
```

In [10]:
```
# Declare an instance of the SVM class.
# Weights are initialized to a random value.
# Note, to keep people's initial solutions consistent, we are going to use
a random seed.

np.random.seed(1)

num_classes = len(np.unique(y_train))
num_features = X_train.shape[1]

svm = SVM(dims=[num_classes, num_features])
```

**SVM loss**

In [11]:
```
## Implement the loss function for in the SVM class(nndl/svm.py), svm.loss(
)

loss = svm.loss(X_train, y_train)
print('The training set loss is {}.'.format(loss))

# If you implemented the loss correctly, it should be 15569.98
```

The training set loss is 15569.9779154.

**SVM gradient**

```
In [12]: ## Calculate the gradient of the SVM class.
         # For convenience, we'll write one function that computes the loss
         #   and gradient together. Please modify svm.loss_and_grad(X, y).
         # You may copy and paste your loss code from svm.loss() here, and then
         #   use the appropriate intermediate values to calculate the gradient.

         loss, grad = svm.loss_and_grad(X_dev,y_dev)

         # Compare your gradient to a numerical gradient check.
         # You should see relative gradient errors on the order of 1e-07 or less if
         you implemented the gradient correctly.
         svm.grad_check_sparse(X_dev, y_dev, grad)
```

```
numerical: -1.206368 analytic: -1.206368, relative error: 2.021268e-08
numerical: -3.718714 analytic: -3.718714, relative error: 4.861115e-08
numerical: -15.580961 analytic: -15.580962, relative error: 8.377393e-09
numerical: 7.447614 analytic: 7.447614, relative error: 9.620118e-09
numerical: 1.463854 analytic: 1.463854, relative error: 2.810277e-08
numerical: -1.582701 analytic: -1.582700, relative error: 5.687247e-08
numerical: 10.859661 analytic: 10.859661, relative error: 2.932411e-08
numerical: -10.477690 analytic: -10.477689, relative error: 8.514592e-09
numerical: 3.381478 analytic: 3.381478, relative error: 1.827548e-09
numerical: -17.365182 analytic: -17.365182, relative error: 2.538693e-09
```

# A vectorized version of SVM

To speed things up, we will vectorize the loss and gradient calculations. This will be helpful for stochastic gradient descent.

```
In [13]: import time
```

In [14]:
```python
## Implement svm.fast_loss_and_grad which calculates the loss and gradient
#    WITHOUT using any for loops.

# Standard loss and gradient
tic = time.time()
loss, grad = svm.loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Normal loss / grad_norm: {} / {} computed in {}s'.format(loss, np.li
nalg.norm(grad, 'fro'), toc - tic))

tic = time.time()
loss_vectorized, grad_vectorized = svm.fast_loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Vectorized loss / grad: {} / {} computed in {}s'.format(loss_vectori
zed, np.linalg.norm(grad_vectorized, 'fro'), toc - tic))

# The losses should match but your vectorized implementation should be much
faster.
print('difference in loss / grad: {} / {}'.format(loss - loss_vectorized, n
p.linalg.norm(grad - grad_vectorized)))

# You should notice a speedup with the same output, i.e., differences on th
e order of 1e-12
```

```
Normal loss / grad_norm: 15879.6858177 / 2182.00689622 computed in 0.074476
9573212s
Vectorized loss / grad: 15879.6858177 / 2182.00689622 computed in 0.0091979
5036316s
difference in loss / grad: 1.27329258248e-11 / 3.77857077752e-12
```
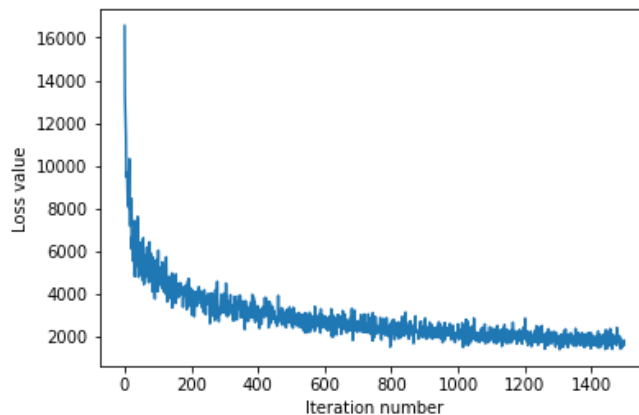
## Stochastic gradient descent

We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

In [15]:
```
# Implement svm.train() by filling in the code to extract a batch of data
# and perform the gradient step.

tic = time.time()
loss_hist = svm.train(X_train, y_train, learning_rate=5e-4,
                      num_iters=1500, verbose=True)
toc = time.time()
print('That took {}s'.format(toc - tic))

plt.plot(loss_hist)
plt.xlabel('Iteration number')
plt.ylabel('Loss value')
plt.show()
```

```
iteration 0 / 1500: loss 16557.3800019
iteration 100 / 1500: loss 4701.08945127
iteration 200 / 1500: loss 4017.33313794
iteration 300 / 1500: loss 3681.9226472
iteration 400 / 1500: loss 2732.6164374
iteration 500 / 1500: loss 2786.63784246
iteration 600 / 1500: loss 2837.03578428
iteration 700 / 1500: loss 2206.23486874
iteration 800 / 1500: loss 2269.03882412
iteration 900 / 1500: loss 2543.23781539
iteration 1000 / 1500: loss 2566.69213573
iteration 1100 / 1500: loss 2182.06890591
iteration 1200 / 1500: loss 1861.11822443
iteration 1300 / 1500: loss 1982.90138585
iteration 1400 / 1500: loss 1927.52041586
That took 5.89598894119s
```



**Evaluate the performance of the trained SVM on the validation data.**

In [16]:
```python
## Implement svm.predict() and use it to compute the training and testing error.

y_train_pred = svm.predict(X_train)
print('training accuracy: {}'.format(np.mean(np.equal(y_train,y_train_pred), )))
y_val_pred = svm.predict(X_val)
print('validation accuracy: {}'.format(np.mean(np.equal(y_val, y_val_pred), ))
```

```
training accuracy: 0.285306122449
validation accuracy: 0.3
```

## Optimize the SVM

Note, to make things faster and simpler, we won't do k-fold cross-validation, but will only optimize the hyperparameters on the validation dataset (X_val, y_val).

In [33]:
```python
# =================================================================== #
# YOUR CODE HERE:
#    Train the SVM with different learning rates and evaluate on the
#      validation data.
#    Report:
#      - The best learning rate of the ones you tested.
#      - The best VALIDATION accuracy corresponding to the best VALIDATION error.
#
#    Select the SVM that achieved the best validation error and report
#      its error rate on the test set.
#    Note: You do not need to modify SVM class for this section
# =================================================================== #

learning_rates = [5e-2, 1e-2, 5e-3, 1e-3, 5e-4, 1e-4, 5e-5, 1e-5]
validation_scores = []


for i in np.arange(len(learning_rates)):
    svm.train(X_train, y_train, learning_rate=learning_rates[i],
                        num_iters=1500, verbose=False)

    y_val_pred = svm.predict(X_val)
    validation_scores.append(np.mean(np.equal(y_val, y_val_pred)))


m_idx = np.argmax(validation_scores)

print validation_scores
print 'Learning rate of {} achieved the best validation rate of {}'.format(
learning_rates[m_idx], validation_scores[m_idx])

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

```
[0.314, 0.275, 0.319, 0.267, 0.303, 0.271, 0.266, 0.214]
Learning rate of 0.005 achieved the best validation rate of 0.319
```

```python
1  import numpy as np
2  import pdb
3
4  """
5  This code was based off of code from cs231n at Stanford University, and modified for ece239as at UCLA.
6  """
7  class SVM(object):
8
9    def __init__(self, dims=[10, 3073]):
10     self.init_weights(dims=dims)
11
12   def init_weights(self, dims):
13     """
14   Initializes the weight matrix of the SVM.  Note that it has shape (C, D)
15   where C is the number of classes and D is the feature size.
16     """
17     self.W = np.random.normal(size=dims)
18
19   def loss(self, X, y):
20     """
21     Calculates the SVM loss.
22
23     Inputs have dimension D, there are C classes, and we operate on minibatches
24     of N examples.
25
26     Inputs:
27     - X: A numpy array of shape (N, D) containing a minibatch of data.
28     - y: A numpy array of shape (N,) containing training labels; y[i] = c means
29       that X[i] has label c, where 0 <= c < C.
30
31     Returns a tuple of:
32     - loss as single float
33     """
34
35     # compute the loss and the gradient
36     num_classes = self.W.shape[0]
37     num_train = X.shape[0]
38     loss = 0.0
39
40     for i in np.arange(num_train):
41
42       # ============================================================= #
43       # YOUR CODE HERE:
44       #   Calculate the normalized SVM loss, and store it as 'loss'.
45       #   (That is, calculate the sum of the losses of all the training
46       #   set margins, and then normalize the loss by the number of
47       #   training examples.)
48       # ============================================================= #
49
50       temp_loss = 0
51       for j in np.arange(num_classes):
52         if j != y[i]:
53           temp_loss += max(0, 1 + self.W[j].dot(X[i]) - self.W[y[i]].dot(X[i]))
54
55       loss += temp_loss
56
57     loss = loss/num_train
58
59       # ============================================================= #
60       # END YOUR CODE HERE
61       # ============================================================= #
62
63     return loss
64
65   def loss_and_grad(self, X, y):
66     """
67   Same as self.loss(X, y), except that it also returns the gradient.
68
69   Output: grad -- a matrix of the same dimensions as W containing
```

```
 70         the gradient of the loss with respect to W.
 71     """
 72
 73     # compute the loss and the gradient
 74     num_classes = self.W.shape[0]
 75     num_train = X.shape[0]
 76     loss = 0.0
 77     grad = np.zeros_like(self.W)
 78
 79     for i in np.arange(num_train):
 80     # ================================================================ #
 81     # YOUR CODE HERE:
 82     #   Calculate the SVM loss and the gradient.  Store the gradient in
 83     #   the variable grad.
 84     # ================================================================ #
 85       temp_loss = 0
 86
 87       for j in np.arange(num_classes):
 88         z = 1 + self.W[j].dot(X[i]) - self.W[y[i]].dot(X[i])
 89
 90         if j != y[i]:
 91           temp_loss += max(0, z)
 92           grad[j] += X[i] if z > 0 else 0
 93           grad[y[i]] += -1*X[i] if z > 0 else 0
 94
 95       loss += temp_loss
 96
 97
 98
 99     # ================================================================ #
100     # END YOUR CODE HERE
101     # ================================================================ #
102
103     loss /= num_train
104     grad /= num_train
105
106     return loss, grad
107
108   def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
109     """
110     sample a few random elements and only return numerical
111     in these dimensions.
112     """
113
114     for i in np.arange(num_checks):
115       ix = tuple([np.random.randint(m) for m in self.W.shape])
116
117       oldval = self.W[ix]
118       self.W[ix] = oldval + h # increment by h
119       fxph = self.loss(X, y)
120       self.W[ix] = oldval - h # decrement by h
121       fxmh = self.loss(X,y) # evaluate f(x - h)
122       self.W[ix] = oldval # reset
123
124       grad_numerical = (fxph - fxmh) / (2 * h)
125       grad_analytic = your_grad[ix]
126       rel_error = abs(grad_numerical - grad_analytic) / (abs(grad_numerical) + abs(grad_analytic))
127       print('numerical: %f analytic: %f, relative error: %e' % (grad_numerical, grad_analytic, rel_error))
128
129   def fast_loss_and_grad(self, X, y):
130     """
131     A vectorized implementation of loss_and_grad. It shares the same
132     inputs and ouputs as loss_and_grad.
133     """
134     loss = 0.0
135     grad = np.zeros(self.W.shape) # initialize the gradient as zero
136
137     # ================================================================ #
138     # YOUR CODE HERE:
139     #   Calculate the SVM loss WITHOUT any for loops.
```

```
140       # ================================================================ #
141
142      loss_mat = self.W.dot(X.T)
143      loss_mat = loss_mat - loss_mat[y,np.arange(loss_mat.shape[1])] + 1
144      loss_mat[y, np.arange(loss_mat.shape[1])] = 0
145      loss_mat[loss_mat < 0] = 0
146      loss = loss_mat.sum()/X.shape[0]
147
148       # ================================================================ #
149       # END YOUR CODE HERE
150       # ================================================================ #
151
152
153       # ================================================================ #
154       # YOUR CODE HERE:
155       #   Calculate the SVM grad WITHOUT any for loops.
156       # ================================================================ #
157      indicator = loss_mat
158      indicator[indicator>0] = 1
159      rsum = np.sum(indicator, axis=0)
160      indicator[y, np.arange(loss_mat.shape[1])] = -rsum
161      grad = indicator.dot(X)/X.shape[0]
162
163       # ================================================================ #
164       # END YOUR CODE HERE
165       # ================================================================ #
166
167      return loss, grad
168
169  def train(self, X, y, learning_rate=1e-3, num_iters=100,
170            batch_size=200, verbose=False):
171      """
172      Train this linear classifier using stochastic gradient descent.
173
174      Inputs:
175      - X: A numpy array of shape (N, D) containing training data; there are N
176        training samples each of dimension D.
177      - y: A numpy array of shape (N,) containing training labels; y[i] = c
178        means that X[i] has label 0 <= c < C for C classes.
179      - learning_rate: (float) learning rate for optimization.
180      - num_iters: (integer) number of steps to take when optimizing
181      - batch_size: (integer) number of training examples to use at each step.
182      - verbose: (boolean) If true, print progress during optimization.
183
184      Outputs:
185      A list containing the value of the loss function at each training iteration.
186      """
187      num_train, dim = X.shape
188      num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is number of classes
189
190      self.init_weights(dims=[np.max(y) + 1, X.shape[1]])  # initializes the weights of self.W
191
192      # Run stochastic gradient descent to optimize W
193      loss_history = []
194
195      for it in np.arange(num_iters):
196        X_batch = None
197        y_batch = None
198
199        # ================================================================ #
200        # YOUR CODE HERE:
201        #   Sample batch_size elements from the training data for use in
202        #   gradient descent.  After sampling,
203        #     - X_batch should have shape: (dim, batch_size)
204        #     - y_batch should have shape: (batch_size,)
205        #   The indices should be randomly generated to reduce correlations
206        #   in the dataset.  Use np.random.choice.  It's okay to sample with
207        #   replacement.
208        # ================================================================ #
209        mask = np.random.choice(np.arange(X.shape[0]), batch_size)
```

```python
210          X_batch = X[mask]
211          y_batch = y[mask]
212          # ================================================================ #
213          # END YOUR CODE HERE
214          # ================================================================ #

216          # evaluate loss and gradient
217          loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
218          loss_history.append(loss)

220          # ================================================================ #
221          # YOUR CODE HERE:
222          #   Update the parameters, self.W, with a gradient step
223          # ================================================================ #

225          self.W -= learning_rate*grad

227          # ================================================================ #
228          # END YOUR CODE HERE
229          # ================================================================ #

231          if verbose and it % 100 == 0:
232            print('iteration {} / {}: loss {}'.format(it, num_iters, loss))

234      return loss_history

236    def predict(self, X):
237      """
238      Inputs:
239      - X: N x D array of training data. Each row is a D-dimensional point.

241      Returns:
242      - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
243        array of length N, and each element is an integer giving the predicted
244        class.
245      """
246      y_pred = np.zeros(X.shape[0])


249      # ================================================================ #
250      # YOUR CODE HERE:
251      #   Predict the labels given the training data with the parameter self.W.
252      # ================================================================ #

254      scores = self.W.dot(X.T)
255      y_pred = np.argmax(scores,axis=0)

257      # ================================================================ #
258      # END YOUR CODE HERE
259      # ================================================================ #

261      return y_pred
```

## This is the softmax workbook for ECE 239AS Assignment #2

Please follow the notebook linearly to implement a softmax classifier.

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and code in the jupyer notebook to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

The goal of this workbook is to give you experience with training a softmax classifier.

```
In [1]:  import random
         import numpy as np
         from cs231n.data_utils import load_CIFAR10
         import matplotlib.pyplot as plt

         %matplotlib inline
         %load_ext autoreload
         %autoreload 2
```

In [2]:
```python
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 prepar
e
    it for the linear classifier. These are the same steps as we used for t
he
    SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = '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,))
```

## Training a softmax classifier.

The following cells will take you through building a softmax classifier. You will implement its loss function, then subsequently train it with gradient descent. Finally, you will choose the learning rate of gradient descent to optimize its classification performance.

```
In [3]:  from nndl import Softmax
```

```
In [4]:  # Declare an instance of the Softmax class.
         # Weights are initialized to a random value.
         # Note, to keep people's first solutions consistent, we are going to use a
         random seed.

         np.random.seed(1)

         num_classes = len(np.unique(y_train))
         num_features = X_train.shape[1]

         softmax = Softmax(dims=[num_classes, num_features])
```

**Softmax loss**

```
In [5]:  ## Implement the loss function of the softmax using a for loop over
         #  the number of examples

         loss = softmax.loss(X_train, y_train)
```

```
In [6]:  print(loss)

         2.3277607028048966
```

## Question:

You'll notice the loss returned by the softmax is about 2.3 (if implemented correctly). Why does this value make sense?

## Answer:

The weights are initialized to a normal distribution with mean 0, which means that the expected value of the samples would most likely also be 0. This would result in the log term in the loss function having an argument of 10 each time, yielding 2.3. Since this should be the same for every sample, the average loss would also come out to 2.3.

**Softmax gradient**

```
In [7]:   ## Calculate the gradient of the softmax loss in the Softmax class.
          # For convenience, we'll write one function that computes the loss
          #   and gradient together, softmax.loss_and_grad(X, y)
          # You may copy and paste your loss code from softmax.loss() here, and then
          #   use the appropriate intermediate values to calculate the gradient.

          loss, grad = softmax.loss_and_grad(X_dev,y_dev)

          # Compare your gradient to a gradient check we wrote.
          # You should see relative gradient errors on the order of 1e-07 or less if
          you implemented the gradient correctly.
          softmax.grad_check_sparse(X_dev, y_dev, grad)
```

```
numerical: -0.461254 analytic: -0.461254, relative error: 1.147387e-09
numerical: 1.415071 analytic: 1.415071, relative error: 3.611131e-08
numerical: 0.221536 analytic: 0.221536, relative error: 1.913397e-08
numerical: 0.823889 analytic: 0.823889, relative error: 3.548639e-08
numerical: 0.131422 analytic: 0.131421, relative error: 5.121652e-07
numerical: 1.079281 analytic: 1.079281, relative error: 4.340823e-08
numerical: -0.442191 analytic: -0.442191, relative error: 1.171090e-07
numerical: -1.603514 analytic: -1.603514, relative error: 1.784821e-08
numerical: -0.774959 analytic: -0.774959, relative error: 1.523634e-08
numerical: -3.820474 analytic: -3.820474, relative error: 1.752965e-08
```

# A vectorized version of Softmax

To speed things up, we will vectorize the loss and gradient calculations. This will be helpful for stochastic gradient descent.

```
In [8]:   import time
```

In [9]:
```python
## Implement softmax.fast_loss_and_grad which calculates the loss and gradi
ent
#    WITHOUT using any for loops.

# Standard loss and gradient
tic = time.time()
loss, grad = softmax.loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Normal loss / grad_norm: {} / {} computed in {}s'.format(loss, np.li
nalg.norm(grad, 'fro'), toc - tic))

tic = time.time()
loss_vectorized, grad_vectorized = softmax.fast_loss_and_grad(X_dev, y_dev)
toc = time.time()
print('Vectorized loss / grad: {} / {} computed in {}s'.format(loss_vectori
zed, np.linalg.norm(grad_vectorized, 'fro'), toc - tic))

# The losses should match but your vectorized implementation should be much
faster.
print('difference in loss / grad: {} /{} '.format(loss - loss_vectorized, n
p.linalg.norm(grad - grad_vectorized)))

# You should notice a speedup with the same output.
```

```
Normal loss / grad_norm: 2.30711514969 / 324.916463415 computed in 0.064689
874649s
Vectorized loss / grad: 2.30711514969 / 324.916463415 computed in 0.0102250
576019s
difference in loss / grad: 4.4408920985e-16 /2.149853952e-13
```

## Stochastic gradient descent

We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

## Question:

How should the softmax gradient descent training step differ from the svm training step, if at all?
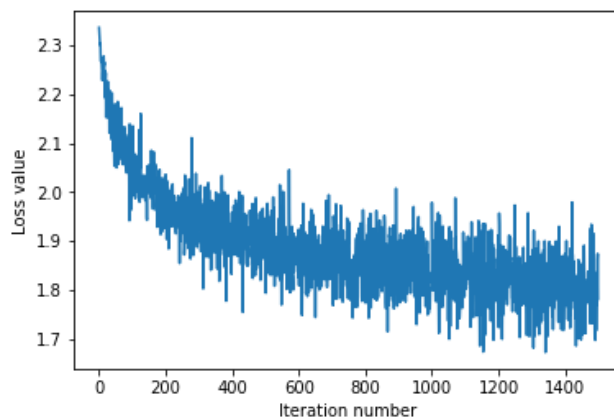
## Answer:

They should be the same, aside from the loss and gradient calculations

```
In [10]:  # Implement softmax.train() by filling in the code to extract a batch of da
          ta
          # and perform the gradient step.
          import time


          tic = time.time()
          loss_hist = softmax.train(X_train, y_train, learning_rate=1e-7,
                                    num_iters=1500, verbose=True)
          toc = time.time()
          print('That took {}s'.format(toc - tic))

          plt.plot(loss_hist)
          plt.xlabel('Iteration number')
          plt.ylabel('Loss value')
          plt.show()
```

```
iteration 0 / 1500: loss 2.33659266066
iteration 100 / 1500: loss 2.05572226139
iteration 200 / 1500: loss 2.03577451207
iteration 300 / 1500: loss 1.98133481656
iteration 400 / 1500: loss 1.9583142444
iteration 500 / 1500: loss 1.86226530735
iteration 600 / 1500: loss 1.85326114544
iteration 700 / 1500: loss 1.83530622237
iteration 800 / 1500: loss 1.82938924688
iteration 900 / 1500: loss 1.89921585304
iteration 1000 / 1500: loss 1.97835035403
iteration 1100 / 1500: loss 1.84707979135
iteration 1200 / 1500: loss 1.84114502687
iteration 1300 / 1500: loss 1.79104024958
iteration 1400 / 1500: loss 1.87058030294
That took 7.12574601173s
```



**Evaluate the performance of the trained softmax classifier on the validation data.**

In [11]: `## Implement softmax.predict() and use it to compute the training and testing error.`

```
y_train_pred = softmax.predict(X_train)
print('training accuracy: {}'.format(np.mean(np.equal(y_train,y_train_pred)
, )))
y_val_pred = softmax.predict(X_val)
print('validation accuracy: {}'.format(np.mean(np.equal(y_val, y_val_pred))
, ))
```

```
training accuracy: 0.381142857143
validation accuracy: 0.398
```

## Optimize the softmax classifier

You may copy and paste your optimization code from the SVM here.

In [12]: `np.finfo(float).eps`

Out[12]: `2.220446049250313e-16`

In [14]:
```
# =================================================================== #
# YOUR CODE HERE:
#    Train the Softmax classifier with different learning rates and
#      evaluate on the validation data.
#    Report:
#      - The best learning rate of the ones you tested.
#      - The best validation accuracy corresponding to the best validation e
rror.
#
#    Select the SVM that achieved the best validation error and report
#      its error rate on the test set.
# =================================================================== #


learning_rates = [5e-5, 1e-5, 5e-6, 1e-6, 5e-7, 1e-7, 5e-8, 1e-8]
validation_scores = []


for i in np.arange(len(learning_rates)):
    softmax.train(X_train, y_train, learning_rate=learning_rates[i],
                  num_iters=1500, verbose=False)

    y_val_pred = softmax.predict(X_val)
    validation_scores.append(np.mean(np.equal(y_val, y_val_pred)))


m_idx = np.argmax(validation_scores)

print validation_scores
print 'Learning rate of {} achieved the best validation rate of {}'.format(
learning_rates[m_idx], validation_scores[m_idx])

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

```
[0.312, 0.346, 0.389, 0.411, 0.404, 0.384, 0.365, 0.313]
Learning rate of 1e-06 achieved the best validation rate of 0.411
```

```python
  1 import numpy as np
  2
  3 class Softmax(object):
  4
  5   def __init__(self, dims=[10, 3073]):
  6     self.init_weights(dims=dims)
  7
  8   def init_weights(self, dims):
  9     """
 10 Initializes the weight matrix of the Softmax classifier.
 11 Note that it has shape (C, D) where C is the number of
 12 classes and D is the feature size.
 13     """
 14     self.W = np.random.normal(size=dims) * 0.0001
 15
 16   def loss(self, X, y):
 17     """
 18     Calculates the softmax loss.
 19
 20     Inputs have dimension D, there are C classes, and we operate on minibatches
 21     of N examples.
 22
 23     Inputs:
 24     - X: A numpy array of shape (N, D) containing a minibatch of data.
 25     - y: A numpy array of shape (N,) containing training labels; y[i] = c means
 26       that X[i] has label c, where 0 <= c < C.
 27
 28     Returns a tuple of:
 29     - loss as single float
 30     """
 31
 32     # Initialize the loss to zero.
 33     loss = 0.0
 34
 35     # ================================================================ #
 36     # YOUR CODE HERE:
 37     #   Calculate the normalized softmax loss.  Store it as the variable loss.
 38     #   (That is, calculate the sum of the losses of all the training
 39     #   set margins, and then normalize the loss by the number of
 40     #   training examples.)
 41     # ================================================================ #
 42
 43     num_samples = X.shape[0]
 44     num_classes = self.W.shape[0]
 45
 46     for i in np.arange(num_samples):
 47
 48       temp_log = 0
 49       for j in np.arange(num_classes):
 50         temp_log += np.exp(self.W[j].dot(X[i]))
 51
 52       loss += np.log(temp_log) - self.W[y[i]].dot(X[i])
 53
 54     loss = loss/num_samples
 55
 56     # ================================================================ #
 57     # END YOUR CODE HERE
 58     # ================================================================ #
 59
 60
 61
 62     return loss
 63
 64   def loss_and_grad(self, X, y):
 65     """
 66 Same as self.loss(X, y), except that it also returns the gradient.
 67
 68 Output: grad -- a matrix of the same dimensions as W containing
 69     the gradient of the loss with respect to W.
```

```python
 70    """
 71
 72        # Initialize the loss and gradient to zero.
 73        loss = 0.0
 74        grad = np.zeros_like(self.W)
 75
 76        # ================================================================ #
 77        # YOUR CODE HERE:
 78        #   Calculate the softmax loss and the gradient. Store the gradient
 79        #   as the variable grad.
 80        # ================================================================ #
 81
 82        num_samples = X.shape[0]
 83        num_classes = self.W.shape[0]
 84
 85        for i in np.arange(num_samples):
 86
 87            temp_log = 0
 88            temp_grad = np.zeros_like(self.W[:,0])
 89            for j in np.arange(num_classes):
 90                score = np.exp(self.W[j].dot(X[i]))
 91                temp_log += score
 92                temp_grad[j] = score
 93
 94            temp_grad /= temp_log
 95            temp_grad[y[i]] -= 1
 96            grad += temp_grad[:, np.newaxis]*X[i]
 97
 98
 99            loss += np.log(temp_log) - self.W[y[i]].dot(X[i])
100
101        loss = loss/num_samples
102        grad = grad/num_samples
103
104
105        # ================================================================ #
106        # END YOUR CODE HERE
107        # ================================================================ #
108
109        return loss, grad
110
111    def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
112        """
113        sample a few random elements and only return numerical
114        in these dimensions.
115        """
116
117        for i in np.arange(num_checks):
118            ix = tuple([np.random.randint(m) for m in self.W.shape])
119
120            oldval = self.W[ix]
121            self.W[ix] = oldval + h # increment by h
122            fxph = self.loss(X, y)
123            self.W[ix] = oldval - h # decrement by h
124            fxmh = self.loss(X,y) # evaluate f(x - h)
125            self.W[ix] = oldval # reset
126
127            grad_numerical = (fxph - fxmh) / (2 * h)
128            grad_analytic = your_grad[ix]
129            rel_error = abs(grad_numerical - grad_analytic) / (abs(grad_numerical) + abs(grad_analytic))
130            print('numerical: %f analytic: %f, relative error: %e' % (grad_numerical, grad_analytic, rel_error))
131
132    def fast_loss_and_grad(self, X, y):
133        """
134        A vectorized implementation of loss_and_grad. It shares the same
135    inputs and ouputs as loss_and_grad.
136        """
137        loss = 0.0
138        grad = np.zeros(self.W.shape) # initialize the gradient as zero
139
```

```python
140        # ================================================================ #
141        # YOUR CODE HERE:
142        #   Calculate the softmax loss and gradient WITHOUT any for loops.
143        # ================================================================ #
144
145        num_samples = X.shape[0]
146        num_classes = self.W.shape[0]
147
148        # Loss
149        scores = self.W.dot(X.T)
150        e_scores = np.exp(self.W.dot(X.T))
151        sums = np.sum(e_scores, axis=0)
152        log_sums = np.log(sums)
153        y_terms = scores[y, np.arange(num_samples)]
154        loss = np.sum(log_sums - y_terms)/num_samples
155
156
157        # Grad
158        e_scores = e_scores/sums
159        e_scores[y,np.arange(num_samples)] -= 1
160        grad = e_scores.dot(X)/num_samples
161
162        # ================================================================ #
163        # END YOUR CODE HERE
164        # ================================================================ #
165
166        return loss, grad
167
168    def train(self, X, y, learning_rate=1e-3, num_iters=100,
169              batch_size=200, verbose=False):
170        """
171        Train this linear classifier using stochastic gradient descent.
172
173        Inputs:
174        - X: A numpy array of shape (N, D) containing training data; there are N
175          training samples each of dimension D.
176        - y: A numpy array of shape (N,) containing training labels; y[i] = c
177          means that X[i] has label 0 <= c < C for C classes.
178        - learning_rate: (float) learning rate for optimization.
179        - num_iters: (integer) number of steps to take when optimizing
180        - batch_size: (integer) number of training examples to use at each step.
181        - verbose: (boolean) If true, print progress during optimization.
182
183        Outputs:
184        A list containing the value of the loss function at each training iteration.
185        """
186        num_train, dim = X.shape
187        num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is number of classes
188
189        self.init_weights(dims=[np.max(y) + 1, X.shape[1]])  # initializes the weights of self.W
190
191        # Run stochastic gradient descent to optimize W
192        loss_history = []
193
194        for it in np.arange(num_iters):
195            X_batch = None
196            y_batch = None
197
198            # ================================================================ #
199            # YOUR CODE HERE:
200            #   Sample batch_size elements from the training data for use in
201            #     gradient descent.  After sampling,
202            #     - X_batch should have shape: (dim, batch_size)
203            #     - y_batch should have shape: (batch_size,)
204            #   The indices should be randomly generated to reduce correlations
205            #   in the dataset.  Use np.random.choice.  It's okay to sample with
206            #   replacement.
207            # ================================================================ #
208            mask = np.random.choice(np.arange(X.shape[0]), batch_size)
209            X_batch = X[mask]
```

```
210         y_batch = y[mask]
211         # =============================================================== #
212         # END YOUR CODE HERE
213         # =============================================================== #
214
215         # evaluate loss and gradient
216         loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
217         loss_history.append(loss)
218
219         # =============================================================== #
220         # YOUR CODE HERE:
221         #   Update the parameters, self.W, with a gradient step
222         # =============================================================== #
223         self.W -= learning_rate*grad
224
225         # =============================================================== #
226         # END YOUR CODE HERE
227         # =============================================================== #
228
229         if verbose and it % 100 == 0:
230           print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
231
232       return loss_history
233
234     def predict(self, X):
235       """
236       Inputs:
237       - X: N x D array of training data. Each row is a D-dimensional point.
238
239       Returns:
240       - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
241         array of length N, and each element is an integer giving the predicted
242         class.
243       """
244       y_pred = np.zeros(X.shape[1])
245       # =============================================================== #
246       # YOUR CODE HERE:
247       #   Predict the labels given the training data.
248       # =============================================================== #
249
250       scores = self.W.dot(X.T)
251       y_pred = np.argmax(scores,axis=0)
252
253
254
255       # =============================================================== #
256       # END YOUR CODE HERE
257       # =============================================================== #
258
259       return y_pred
```