

```

import numpy as np
import pdb

class KNN(object):

    def __init__(self):
        pass

    def train(self, X, y):
        """
        Inputs:
        - X is a numpy array of size (num_examples, D)
        - y is a numpy array of size (num_examples, )
        """
        self.X_train = X
        self.y_train = y

    def compute_distances(self, X, norm=None):
        """
        Compute the distance between each test point in X and each training point
        in self.X_train.

        Inputs:
        - X: A numpy array of shape (num_test, D) containing test data.
        - norm: the function with which the norm is taken.

        Returns:
        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
          is the Euclidean distance between the ith test point and the jth training
          point.
        """
        if norm is None:
            norm = lambda x: np.sqrt(np.sum(x**2))
            #norm = 2

        num_test = X.shape[0]
        num_train = self.X_train.shape[0]
        dists = np.zeros((num_test, num_train))
        for i in np.arange(num_test):

            for j in np.arange(num_train):
                # ===== #
                # YOUR CODE HERE:
                #   Compute the distance between the ith test point and the jth
                #   training point using norm(), and store the result in dists[i, j].
                # ===== #

                trainSample=self.X_train[j]
                testSample = X[i]
                distance = norm(trainSample-testSample)
                dists[i, j] = distance

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

```

```

    return dists

def compute_L2_distances_vectorized(self, X):
    """
    Compute the distance between each test point in X and each training point
    in self.X_train WITHOUT using any for loops.

    Inputs:
    - X: A numpy array of shape (num_test, D) containing test data.

    Returns:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      is the Euclidean distance between the ith test point and the jth training
      point.
    """
    num_test = X.shape[0]
    num_train = self.X_train.shape[0]
    dists = np.zeros((num_test, num_train))

    # ===== #
    # YOUR CODE HERE:
    # Compute the L2 distance between the ith test point and the jth
    # training point and store the result in dists[i, j]. You may
    # NOT use a for loop (or list comprehension). You may only use
    # numpy operations.
    #
    # HINT: use broadcasting. If you have a shape (N,1) array and
    # a shape (M,) array, adding them together produces a shape (N, M)
    # array.
    # ===== #

    trainsquared = np.sum(self.X_train**2, axis=1, keepdims=True)
    testsquared = np.sum(X**2, axis=1)
    multiplied = np.dot(self.X_train, X.T)
    dists = np.sqrt(trainsquared - 2*multiplied + testsquared)
    dists = dists.T

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

    return dists

def compute_L1_distances_vectorized(self, X):
    """
    Compute the distance between each test point in X and each training point
    in self.X_train WITHOUT using any for loops.

    Inputs:
    - X: A numpy array of shape (num_test, D) containing test data.

    Returns:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      is the Euclidean distance between the ith test point and the jth training
      point.
    """
    num_test = X.shape[0]

```

```

num_train = self.X_train.shape[0]
dists = np.zeros((num_test, num_train))

# ===== #
# YOUR CODE HERE:
#   Compute the L2 distance between the ith test point and the jth
#   training point and store the result in dists[i, j]. You may
#   NOT use a for loop (or list comprehension). You may only use
#   numpy operations.
#
#   HINT: use broadcasting. If you have a shape (N,1) array and
#   a shape (M,) array, adding them together produces a shape (N, M)
#   array.
# ===== #
for i in range(0,num_test):
    # find the nearest training image to the i'th test image
    # using the L1 distance (sum of absolute value differences)
    distances = np.sum(np.abs(self.X_train - X[i,:]), axis = 1)
    dists[i]=distances

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

return dists

def compute_Linf_distances_vectorized(self, X):
    """
    Compute the distance between each test point in X and each training point
    in self.X_train WITHOUT using any for loops.

    Inputs:
    - X: A numpy array of shape (num_test, D) containing test data.

    Returns:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      is the Euclidean distance between the ith test point and the jth training
      point.
    """
    num_test = X.shape[0]
    num_train = self.X_train.shape[0]
    dists = np.zeros((num_test, num_train))

    # ===== #
    # YOUR CODE HERE:
    #   Compute the L2 distance between the ith test point and the jth
    #   training point and store the result in dists[i, j]. You may
    #   NOT use a for loop (or list comprehension). You may only use
    #   numpy operations.
    #
    #   HINT: use broadcasting. If you have a shape (N,1) array and
    #   a shape (M,) array, adding them together produces a shape (N, M)
    #   array.
    # ===== #
    for i in range(0,num_test):
        # find the nearest training image to the i'th test image

```

```

# using the L1 distance (sum of absolute value differences)
distances = np.sum(np.linalg.norm(self.X_train - X[i,:],ord="inf"), axis = 1)
dists[i]=distances

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

return dists

def predict_labels(self, dists, k=1):
    """
    Given a matrix of distances between test points and training points,
    predict a label for each test point.

    Inputs:
    - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
      gives the distance between the ith test point and the jth training point.

    Returns:
    - y: A numpy array of shape (num_test,) containing predicted labels for the
      test data, where y[i] is the predicted label for the test point X[i].
    """
    num_test = dists.shape[0]
    y_pred = np.zeros(num_test)
    for i in np.arange(num_test):
        # A list of length k storing the labels of the k nearest neighbors to
        # the ith test point.
        closest_y = []
        # ===== #
        # YOUR CODE HERE:
        # Use the distances to calculate and then store the labels of
        # the k-nearest neighbors to the ith test point. The function
        # numpy.argsort may be useful.
        #
        # After doing this, find the most common label of the k-nearest
        # neighbors. Store the predicted label of the ith training example
        # as y_pred[i]. Break ties by choosing the smaller label.
        # ===== #
        distance = dists[i]
        sortIndices = np.argsort(distance)
        best0nes = sortIndices[0:k]
        closest_y = self.y_train[best0nes]
        counts = np.bincount(closest_y)
        y_pred[i]=np.argmax(counts)
        # ===== #
        # END YOUR CODE HERE
        # ===== #

    return y_pred

```

# This is the k-nearest neighbors workbook for ECE C147/C247 Assignment #2

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

Please print out the workbook entirely when completed.

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]: 1 import numpy as np # for doing most of our calculations
2 import matplotlib.pyplot as plt # for plotting
3 from utils.data_utils import load_CIFAR10 # function to load the CIFAR-10 data
4
5 # Load matplotlib images inline
6 %matplotlib inline
7
8 # These are important for reloading any code you write in external .py files
9 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
10 %load_ext autoreload
11 %autoreload 2
```

```
In [2]: 1 # Set the path to the CIFAR-10 data
2 cifar10_dir = 'dataset\cifar-10-batches-py' # You need to update this line
3 X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
4
5 # As a sanity check, we print out the size of the training and test data.
6 print('Training data shape: ', X_train.shape)
7 print('Training labels shape: ', y_train.shape)
8 print('Test data shape: ', X_test.shape)
9 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]:

```

1  # Visualize some examples from the dataset.
2  # We show a few examples of training images from each class.
3  classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 's
4  num_classes = len(classes)
5  samples_per_class = 7
6  for y, cls in enumerate(classes):
7      idxs = np.flatnonzero(y_train == y)
8      idxs = np.random.choice(idxs, samples_per_class, replace=False)
9      for i, idx in enumerate(idxs):
10         plt_idx = i * num_classes + y + 1
11         plt.subplot(samples_per_class, num_classes, plt_idx)
12         plt.imshow(X_train[idx].astype('uint8'))
13         plt.axis('off')
14         if i == 0:
15             plt.title(cls)
16 plt.show()

```



```
In [4]: 1 # Subsample the data for more efficient code execution in this exercise
2 num_training = 5000
3 mask = list(range(num_training))
4 X_train = X_train[mask]
5 y_train = y_train[mask]
6
7 num_test = 500
8 mask = list(range(num_test))
9 X_test = X_test[mask]
10 y_test = y_test[mask]
11
12 # Reshape the image data into rows
13 X_train = np.reshape(X_train, (X_train.shape[0], -1))
14 X_test = np.reshape(X_test, (X_test.shape[0], -1))
15 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]: 1 # Import the KNN class
2
3 from nnd1 import KNN
```

```
In [6]: 1 # Declare an instance of the knn class.
2 knn = KNN()
3
4 # Train the classifier.
5 # We have implemented the training of the KNN classifier.
6 # Look at the train function in the KNN class to see what this does.
7 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) In the `knn.train()` function, we are simply storing all the training data samples and the corresponding labels in memory.
- (2) An obvious pro is that the training step is extremely fast and simple to implement. However, a con is memory usage. As the amount of samples increase, the model needs more storage.

## 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]: 1 # Implement the function compute_distances() in the KNN class.
2 # Do not worry about the input 'norm' for now; use the default definition of
3 # in the code, which is the 2-norm.
4 # You should only have to fill out the clearly marked sections.
5
6 import time
7 time_start =time.time()
8
9 dists_L2 = knn.compute_distances(X=X_test)
10
11 print('Time to run code: {}'.format(time.time()-time_start))
12 print('Frobenius norm of L2 distances: {}'.format(np.linalg.norm(dists_L2, 'fro')))
```

Time to run code: 47.018675088882446

Frobenius norm of L2 distances: 7906696.077040902

### 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]: 1 # Implement the function compute_L2_distances_vectorized() in the KNN class.
2 # In this function, you ought to achieve the same L2 distance but WITHOUT an
3 # Note, this is SPECIFIC for the L2 norm.
4
5 time_start =time.time()
6 dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_test)
7 print('Time to run code: {}'.format(time.time()-time_start))
8 print('Difference in L2 distances between your KNN implementations (should be 0): 0.0')
```

Time to run code: 0.5874133110046387

Difference in L2 distances between your KNN implementations (should be 0): 0.0

```
In [12]: 1 print(dists_L2_vectorized.shape)
```

(500, 5000)

### 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]: 1 # Implement the function predict_labels in the KNN class.
2 # Calculate the training error (num_incorrect / total_samples)
3 # from running knn.predict_labels with k=1
4
5 error = 1
6
7 # ===== #
8 # YOUR CODE HERE:
9 # Calculate the error rate by calling predict_labels on the test
10 # data with k = 1. Store the error rate in the variable error.
11 # ===== #
12 distances = knn.compute_L2_distances_vectorized(X=X_test)
13 predictions = knn.predict_labels(distances,1)
14 error = sum(predictions!=y_test)/len(y_test)
15 # ===== #
16 # END YOUR CODE HERE
17 # ===== #
18
19 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]: 1 # Create the dataset folds for cross-validation.
2 num_folds = 5
3
4 X_train_folds = []
5 y_train_folds = []
6 np.random.seed(24)
7 # ===== #
8 # YOUR CODE HERE:
9 # Split the training data into num_folds (i.e., 5) folds.
10 # X_train_folds is a list, where X_train_folds[i] contains the
11 # data points in fold i.
12 # y_train_folds is also a list, where y_train_folds[i] contains
13 # the corresponding labels for the data in X_train_folds[i]
14 # ===== #
15 indices = np.arange(0,X_train.shape[0])
16 randomIndices=np.random.permutation(indices)
17 foldSize = int(X_train.shape[0]/num_folds)
18 start = 0
19 end = start + foldSize
20 for i in range(0,num_folds-1):
21     indices = randomIndices[start:end]
22     foldtraindata=X_train[indices]
23     foldtrainlabel = y_train[indices]
24     X_train_folds.append(foldtraindata)
25     y_train_folds.append(foldtrainlabel)
26     start = end
27     end = start+foldSize
28 indices = randomIndices[start:]
29 foldtraindata=X_train[indices]
30 foldtrainlabel = y_train[indices]
31 X_train_folds.append(foldtraindata)
32 y_train_folds.append(foldtrainlabel)
33
34 # ===== #
35 # END YOUR CODE HERE
36 # ===== #
37
38

```

## 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.

```

In [11]: 1 time_start =time.time()
2
3 ks = [1, 2, 3, 5, 7, 10, 15, 20, 25, 30]
4
5 # ===== #
6 # YOUR CODE HERE:
7 # Calculate the cross-validation error for each k in ks, testing
8 # the trained model on each of the 5 folds. Average these errors
9 # together and make a plot of k vs. cross-validation error. Since
10 # we are assuming L2 distance here, please use the vectorized code!
11 # Otherwise, you might be waiting a long time.
12 # ===== #
13 knn = KNN()
14 kError = []
15 for k in ks:
16     foldError = []
17     for i in range(0,len(X_train_folds)):
18         trainCopy = X_train_folds.copy()
19         trainLabelsCopy = y_train_folds.copy()
20         foldTest = trainCopy[i]
21         foldTestlabel = trainLabelsCopy[i]
22         trainCopy.pop(i)
23         trainLabelsCopy.pop(i)
24         foldTr=np.array(trainCopy)
25         foldTrLab=np.array(trainLabelsCopy)
26         foldTr = foldTr.reshape([foldTr.shape[0]*foldTr.shape[1],foldTr.shap
27         foldTrLab = foldTrLab.reshape([foldTrLab.shape[0]*foldTrLab.shape[1]
28         knn.train(foldTr,foldTrLab)
29         distances = knn.compute_L2_distances_vectorized(X=foldTest)
30         predictions = knn.predict_labels(distances,k)
31         error = sum(predictions!=foldTestlabel)/len(foldTestlabel)
32         foldError.append(error)
33     print("Done cross validation for k %d. Error is %.2f"%(k,np.mean(foldErr
34     kError.append(np.mean(foldError))
35
36
37 plt.plot(ks, kError, 'bo--')
38 plt.axis([0, 31, 0.71, 0.77])
39 plt.xlabel('k value')
40 plt.ylabel('mean cross validation error')
41 plt.show()
42 # ===== #
43 # END YOUR CODE HERE
44 # ===== #
45
46 print('Computation time: %.2f'%(time.time()-time_start))

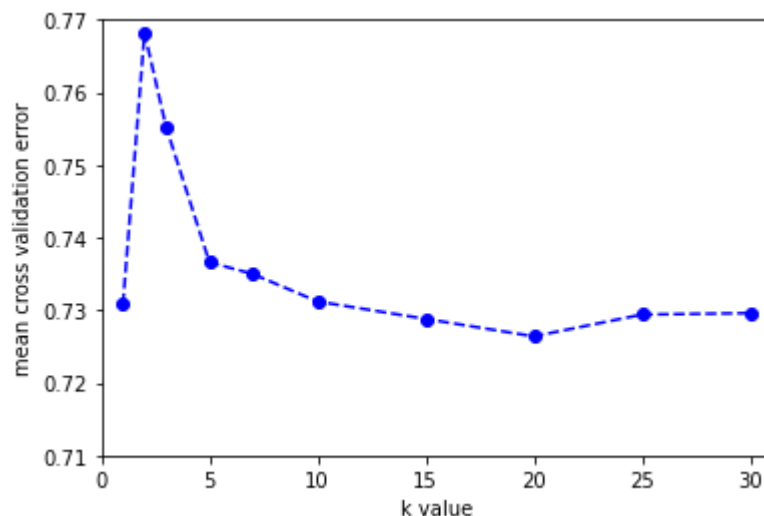
```

```

Done cross validation for k 1. Error is 0.73
Done cross validation for k 2. Error is 0.77
Done cross validation for k 3. Error is 0.76
Done cross validation for k 5. Error is 0.74
Done cross validation for k 7. Error is 0.73
Done cross validation for k 10. Error is 0.73
Done cross validation for k 15. Error is 0.73
Done cross validation for k 20. Error is 0.73

```

Done cross validation for k 25. Error is 0.73  
Done cross validation for k 30. Error is 0.73



Computation time: 35.01

In [12]:

1	kError
---	--------

Out[12]: [0.7307999999999999,  
0.7681999999999999,  
0.7552000000000001,  
0.7365999999999999,  
0.735,  
0.7312,  
0.7288,  
0.7264,  
0.7293999999999999,  
0.7295999999999999]

## 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) The best value of  $k$  is found to be 20.

(2) The cross validation error is 0.7264.

## 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 [13]: 1 time_start =time.time()
2
3 L1_norm = lambda x: np.linalg.norm(x, ord=1)
4 L2_norm = lambda x: np.linalg.norm(x, ord=2)
5 Linf_norm = lambda x: np.linalg.norm(x, ord= np.inf)
6 norms = [L1_norm, L2_norm, Linf_norm]
7 normNames=["l1", "l2", "inf"]
8 # ===== #
9 # YOUR CODE HERE:
10 # Calculate the cross-validation error for each norm in norms, testing
11 # the trained model on each of the 5 folds. Average these errors
12 # together and make a plot of the norm used vs the cross-validation error
13 # Use the best cross-validation k from the previous part.
14 #
15 # Feel free to use the compute_distances function. We're testing just
16 # three norms, but be advised that this could still take some time.
17 # You're welcome to write a vectorized form of the L1- and Linf- norms
18 # to speed this up, but it is not necessary.
19 # ===== #
20 bestK=20
21 knn = KNN()
22 normError = []
23 for ind in range(0,len(norms)):
24     normType = norms[ind]
25     normName = normNames[ind]
26     foldError = []
27     for i in range(0,len(X_train_folds)):
28         trainCopy = X_train_folds.copy()
29         trainLabelsCopy = y_train_folds.copy()
30         foldTest = trainCopy[i]
31         foldTtestlabel = trainLabelsCopy[i]
32         trainCopy.pop(i)
33         trainLabelsCopy.pop(i)
34         foldTr=np.array(trainCopy)
35         foldTrLab=np.array(trainLabelsCopy)
36         foldTr = foldTr.reshape([foldTr.shape[0]*foldTr.shape[1],foldTr.shap
37         foldTrLab = foldTrLab.reshape([foldTrLab.shape[0]*foldTrLab.shape[1]
38         knn.train(foldTr,foldTrLab)
39         if normName=="l2":
40             distances = knn.compute_L2_distances_vectorized(foldTest)
41         else:
42             distances = knn.compute_distances(foldTest,normType)
43
44         #distances = knn.compute_distances(X=foldTest,norm=normType)
45         predictions = knn.predict_labels(distances,k)
46         error = sum(predictions!=foldTtestlabel)/len(foldTtestlabel)
47         print("Fold done")
48         foldError.append(error)
49     print("Done cross validation for norm %. Error is %.2f"%(normName,np.me
50     normError.append(np.mean(foldError))
51 plt.figure()
52 normNames=['L1_norm', 'L2_norm', 'Linf_norm']
53 plt.plot(normNames,normError)
54 plt.xlabel('Norm Type')
55 plt.ylabel('Mean Cross Val Error')
56 # ===== #

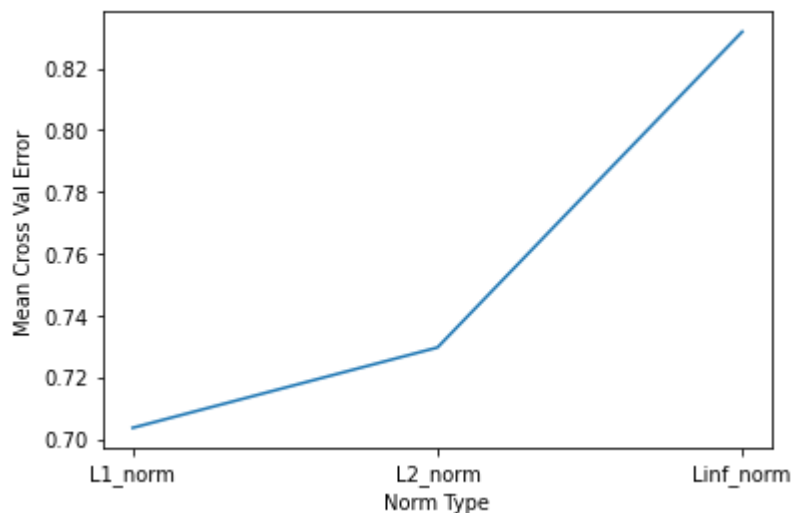
```

```

57 # END YOUR CODE HERE
58 # ===== #
59 print('Computation time: %.2f'%(time.time()-time_start))

```

Fold done  
 Fold done  
 Fold done  
 Fold done  
 Fold done  
 Done cross validation for norm l1. Error is 0.70  
 Fold done  
 Fold done  
 Fold done  
 Fold done  
 Fold done  
 Done cross validation for norm l2. Error is 0.73  
 Fold done  
 Fold done  
 Fold done  
 Fold done  
 Done cross validation for norm inf. Error is 0.83  
 Computation time: 711.44



In [64]: 1 normError

Out[64]: [0.7036, 0.7295999999999999, 0.8318]

## 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) L1 norm has the best cross-validation error.

(2) It is 0.7036

## 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 [14]: 1 error = 1
          2
          3 # ===== #
          4 # YOUR CODE HERE:
          5 #   Evaluate the testing error of the k-nearest neighbors classifier
          6 #   for your optimal hyperparameters found by 5-fold cross-validation.
          7 # ===== #
          8 bestK=20
          9 bestKnn=KNN()
         10 bestKnn.train(X=X_train, y=y_train)
         11 distances = bestKnn.compute_distances(X=X_test,norm=L1_norm)
         12 predictions = bestKnn.predict_labels(distances,bestK)
         13 error = sum(predictions!=y_test)/len(y_test)
         14
         15 # ===== #
         16 # END YOUR CODE HERE
         17 # ===== #
         18
         19 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:

The naive model with  $k=1$  had a test error of 0.726. By implementing cross validation and choosing the best  $k$  and norm values, we managed to reduce the error to 0.72. Therefore, we improved the performance by 0.06.



## ECE C247 HW #2

### 2) Softmax Classifier Gradient

Let  $w_i^T x^{(j)} + b_i = 0_i$

$$\text{Likelihood} = \prod_{j=1}^m \prod_{n=1}^c P_r(y^{(j)} = n | x^{(j)}, \theta) \delta_n(y^{(j)})$$

where  $\delta_n(y^{(j)}) = 1$  if  $y^{(j)} = n$   
0 otherwise.

We are given that  $P_r(y^{(j)} = n | x^{(j)}, \theta) = \text{softmax}_n(x^{(j)}) = \frac{e^{\theta_n}}{\sum_{k=0}^c e^{\theta_k}}$

$$\text{So, likelihood} = \prod_{j=1}^m \prod_{n=1}^c \text{softmax}_n(x^{(j)}) \delta_n(y^{(j)})$$

$$\mathcal{L} = \text{Log-likelihood} = \sum_{j=1}^m \sum_{n=1}^c \delta_n(y^{(j)}) \log(\text{softmax}_n(x^{(j)}))$$

$$\frac{d\mathcal{L}}{dw_i} = \frac{d\mathcal{L}}{d\text{softmax}_n} \cdot \frac{d\text{softmax}_n}{d\theta_i} \cdot \frac{d\theta_i}{dw_i}$$



$$\frac{dL}{do_i} = \sum_{j=1}^m \delta_i(y^{(j)}) \cdot \frac{1}{\text{Softmax}_i(x^{(j)})} \cdot \frac{d \text{Softmax}_i(x^{(j)})}{do_i} + \sum_{n \neq i} \delta_n(y^{(j)}) \cdot \frac{1}{\text{Softmax}_n(x^{(j)})} \cdot \frac{d \text{Softmax}_n(x^{(j)})}{do_i}$$

$$\begin{aligned} \frac{d \text{Softmax}_i(x^{(j)})}{do_i} &= \frac{d e^{o_i}}{\sum_{k=1}^c e^{o_k}} = \frac{e^{o_i} \cdot \sum_{k=1}^c e^{o_k} - e^{o_i} \cdot e^{o_i}}{\left( \sum_{k=1}^c e^{o_k} \right)^2} \\ &= \frac{e^{o_i}}{\sum_k e^{o_k}} - \left( \frac{e^{o_i}}{\sum_k e^{o_k}} \right)^2 \end{aligned}$$

$$= \text{Softmax}_i(x^{(j)}) [1 - \text{Softmax}_i(x^{(j)})]$$

$$\frac{d \text{Softmax}_n(x^{(j)})}{do_i} \text{ when } n \neq i = \frac{d e^{o_n}}{\sum_k e^{o_k}}$$

$$= \frac{0 \cdot \sum_k e^{o_k} - e^{o_n} \cdot e^{o_i}}{\left( \sum_k e^{o_k} \right)^2}$$

$$= -\text{Softmax}_n(x^{(j)}) \cdot \text{Softmax}_i(x^{(j)})$$

Plugging in to  $\frac{dL}{do_i}$  gives:



$$\frac{dL}{do_i} = \sum_{j=1}^m f_i(y^{(j)}) \cdot \frac{\text{softmax}_i(x^{(j)}) \cdot [1 - \text{softmax}_i(x^{(j)})]}{\text{softmax}_i(x^{(j)})} + \sum_{n \neq i} f_n(y^{(j)}) \cdot \frac{-\text{softmax}_n(x^{(j)}) \cdot \text{softmax}_i(x^{(j)})}{\text{softmax}_n(x^{(j)})}$$

$$= \sum_{j=1}^m f_i(y^{(j)}) \cdot [1 - \text{softmax}_i(x^{(j)})] + \sum_{n \neq i} f_n(y^{(j)}) \cdot -\text{softmax}_i(x^{(j)})$$

$$= \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)}) \cdot \sum_{n=1}^n f_n(y^{(j)})$$

We know that  $\sum_{n=1}^n f_n(y^{(j)}) = 1$  since  $f_n(y^{(j)})$  is 1 when  $y^{(j)} = n$  0 otherwise.

$$\text{So, } \frac{dL}{do_i} = \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})$$

$$\frac{do_i}{dw_i} = \frac{dw_i^T x^{(j)} + b_i}{dw_i} = x^{(j)}$$

$$\frac{do_i}{db_i} = \frac{dw_i^T x^{(j)} + b_i}{db_i} = 1$$

$$\text{So, } \frac{dL}{dw_i} = \frac{dL}{do_i} \cdot \frac{do_i}{dw_i} = \sum_{j=1}^m [f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})] \cdot x^{(j)}$$

$$\frac{dL}{db_i} = \frac{dL}{do_i} \cdot \frac{do_i}{db_i} = \sum_{j=1}^m f_i(y^{(j)}) - \text{softmax}_i(x^{(j)})$$

```

import numpy as np

class Softmax(object):

    def __init__(self, dims=[10, 3073]):
        self.init_weights(dims=dims)

    def init_weights(self, dims):
        """
        Initializes the weight matrix of the Softmax classifier.
        Note that it has shape (C, D) where C is the number of
        classes and D is the feature size.
        """
        self.W = np.random.normal(size=dims) * 0.0001

    def loss(self, X, y):
        """
        Calculates the softmax loss.

        Inputs have dimension D, there are C classes, and we operate on minibatches
        of N examples.

        Inputs:
        - X: A numpy array of shape (N, D) containing a minibatch of data.
        - y: A numpy array of shape (N,) containing training labels; y[i] = c means
            that X[i] has label c, where 0 ≤ c < C.

        Returns a tuple of:
        - loss as single float
        """
        # Initialize the loss to zero.
        loss = 0.0

        # ===== #
        # YOUR CODE HERE:
        # Calculate the normalized softmax loss. Store it as the variable loss.
        # (That is, calculate the sum of the losses of all the training
        # set margins, and then normalize the loss by the number of
        # training examples.)
        # ===== #
        weights = self.W
        losses = []

        for index in range(0, X.shape[0]):
            data = X[index]
            multiplication = weights.dot(data.T).T
            exponentials = np.exp(multiplication)
            exponentialSum = np.sum(exponentials)
            exponentials = exponentials / exponentialSum
            label = y[index]
            crossEntropy = -np.log(exponentials[label])
            losses.append(crossEntropy)
        loss = sum(losses) / len(losses)

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

```

```

# ===== #

return loss

def loss_and_grad(self, X, y):
    """
    Same as self.loss(X, y), except that it also returns the gradient.

    Output: grad -- a matrix of the same dimensions as W containing
    the gradient of the loss with respect to W.
    """

    # Initialize the loss and gradient to zero.
    loss = 0.0
    grad = np.zeros_like(self.W)

    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and the gradient. Store the gradient
    # as the variable grad.
    # ===== #
    weights = self.W
    losses=[]
    for index in range(0,X.shape[0]):
        data = X[index]
        multiplication = weights.dot(data.T).T
        exponentials = np.exp(multiplication)
        exponentialSum = np.sum(exponentials)
        exponentials=exponentials/exponentialSum
        label = y[index]
        crossEntropy = -np.log(exponentials[label])
        losses.append(crossEntropy)
        yHot = np.zeros([1,np.max(y)+1])
        yHot[0,label]=1
        data = np.reshape(data,[1,data.shape[0]])
        exponentials = np.reshape(exponentials,[1,exponentials.shape[0]])
        derivative = np.dot(data.T, (exponentials - yHot))
        grad=grad+derivative.T
    loss = sum(losses)/len(losses)
    grad = grad/[X.shape[0]]

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

    return loss, grad

def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
    """
    sample a few random elements and only return numerical
    in these dimensions.
    """

    for i in np.arange(num_checks):
        ix = tuple([np.random.randint(m) for m in self.W.shape])

        oldval = self.W[ix]

```

```

self.W[ix] = oldval + h # increment by h
fxph = self.loss(X, y)
self.W[ix] = oldval - h # decrement by h
fxmh = self.loss(X, y) # evaluate f(x - h)
self.W[ix] = oldval # reset

grad_numerical = (fxph - fxmh) / (2 * h)
grad_analytic = your_grad[ix]
rel_error = abs(grad_numerical - grad_analytic) / (abs(grad_numerical) + abs(grad_analytic))
print('numerical: %f analytic: %f, relative error: %e' % (grad_numerical, grad_analytic, rel_error))

def fast_loss_and_grad(self, X, y):
    """
    A vectorized implementation of loss_and_grad. It shares the same
    inputs and outputs as loss_and_grad.
    """
    loss = 0.0
    grad = np.zeros(self.W.shape) # initialize the gradient as zero

    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and gradient WITHOUT any for loops.
    # ===== #
    weights = self.W
    multiplication = weights.dot(X.T).T
    multiplication = (multiplication.T - np.amax(multiplication, axis = 1)).T
    soft = np.exp(multiplication)
    sums = np.sum(soft, axis=1)
    probs = (soft.T / sums).T
    predsForClass = probs[np.arange(y.size), y]
    loss = -np.log(predsForClass + 1e-10) # To avoid log(0)
    loss = np.mean(loss)
    yHot = np.zeros([y.shape[0], 9+1])
    yHot[np.arange(y.size), y] = 1
    grad = (1/X.shape[0])*np.dot(X.T, (probs - yHot)).T

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

    return loss, grad

def train(self, X, y, learning_rate=1e-3, num_iters=100,
          batch_size=200, verbose=False):
    """
    Train this linear classifier using stochastic gradient descent.

    Inputs:
    - X: A numpy array of shape (N, D) containing training data; there are N
        training samples each of dimension D.
    - y: A numpy array of shape (N,) containing training labels; y[i] = c
        means that X[i] has label 0 <= c < C for C classes.
    - learning_rate: (float) learning rate for optimization.
    - num_iters: (integer) number of steps to take when optimizing
    - batch_size: (integer) number of training examples to use at each step.
    - verbose: (boolean) If true, print progress during optimization.
    """

```

Outputs:

A list containing the value of the loss function at each training iteration.

"""

num\_train, dim = X.shape

num\_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is number of classes

self.init\_weights(dims=[np.max(y) + 1, X.shape[1]]) # initializes the weights of self.W

# Run stochastic gradient descent to optimize W

loss\_history = []

for it in np.arange(num\_iters):

    X\_batch = None

    y\_batch = None

    # ===== #

    # YOUR CODE HERE:

    # Sample batch\_size elements from the training data for use in

    # gradient descent. After sampling,

    # - X\_batch should have shape: (dim, batch\_size)

    # - y\_batch should have shape: (batch\_size,)

    # The indices should be randomly generated to reduce correlations

    # in the dataset. Use np.random.choice. It's okay to sample with

    # replacement.

    # ===== #

    indices = np.random.choice(np.arange(num\_train), batch\_size)

    XBatch = X[indices]

    yBatch = y[indices]

    # ===== #

    # END YOUR CODE HERE

    # ===== #

    # evaluate loss and gradient

    loss, grad = self.fast\_loss\_and\_grad(XBatch, yBatch)

    loss\_history.append(loss)

    # ===== #

    # YOUR CODE HERE:

    # Update the parameters, self.W, with a gradient step

    # ===== #

    self.W = self.W - learning\_rate\*grad

    # ===== #

    # END YOUR CODE HERE

    # ===== #

    if verbose and it % 100 == 0:

        print('iteration {} / {}: loss {}'.format(it, num\_iters, loss))

return loss\_history

def predict(self, X):

"""

Inputs:

- X: N x D array of training data. Each row is a D-dimensional point.

Returns:

- `y_pred`: Predicted labels for the data in `X`. `y_pred` is a 1-dimensional array of length `N`, and each element is an integer giving the predicted class.

```
"""
y_pred = np.zeros(X.shape[1])
# ===== #
# YOUR CODE HERE:
#   Predict the labels given the training data.
# ===== #
weights = self.W
multiplication = weights.dot(X.T).T
soft = np.exp(multiplication)
y_pred = np.argmax(soft, axis=1)
# ===== #
# END YOUR CODE HERE
# ===== #

return y_pred
```



## This is the softmax workbook for ECE C147/C247 Assignment #2

Please follow the notebook linearly to implement a softmax classifier.

Please print out the workbook entirely when completed.

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

In [1]:

```
1 import random
2 import numpy as np
3 from utils.data_utils import load_CIFAR10
4 import matplotlib.pyplot as plt
5
6 %matplotlib inline
7 %load_ext autoreload
8 %autoreload 2
```

```

In [5]: 1 def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000,
2         """
3         Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
4         it for the linear classifier. These are the same steps as we used for th
5         SVM, but condensed to a single function.
6         """
7         # Load the raw CIFAR-10 data
8         cifar10_dir = 'dataset\cifar-10-batches-py' # You need to update this li
9         X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
10
11        # subsample the data
12        mask = list(range(num_training, num_training + num_validation))
13        X_val = X_train[mask]
14        y_val = y_train[mask]
15        mask = list(range(num_training))
16        X_train = X_train[mask]
17        y_train = y_train[mask]
18        mask = list(range(num_test))
19        X_test = X_test[mask]
20        y_test = y_test[mask]
21        mask = np.random.choice(num_training, num_dev, replace=False)
22        X_dev = X_train[mask]
23        y_dev = y_train[mask]
24
25        # Preprocessing: reshape the image data into rows
26        X_train = np.reshape(X_train, (X_train.shape[0], -1))
27        X_val = np.reshape(X_val, (X_val.shape[0], -1))
28        X_test = np.reshape(X_test, (X_test.shape[0], -1))
29        X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))
30
31        # Normalize the data: subtract the mean image
32        mean_image = np.mean(X_train, axis = 0)
33        X_train -= mean_image
34        X_val -= mean_image
35        X_test -= mean_image
36        X_dev -= mean_image
37
38        # add bias dimension and transform into columns
39        X_train = np.hstack([X_train, np.ones((X_train.shape[0], 1))])
40        X_val = np.hstack([X_val, np.ones((X_val.shape[0], 1))])
41        X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
42        X_dev = np.hstack([X_dev, np.ones((X_dev.shape[0], 1))])
43
44        return X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev
45
46
47        # Invoke the above function to get our data.
48        X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev = get_CIFAR10_d
49        print('Train data shape: ', X_train.shape)
50        print('Train labels shape: ', y_train.shape)
51        print('Validation data shape: ', X_val.shape)
52        print('Validation labels shape: ', y_val.shape)
53        print('Test data shape: ', X_test.shape)
54        print('Test labels shape: ', y_test.shape)
55        print('dev data shape: ', X_dev.shape)
56        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 [6]: 1 from nndl import Softmax
```

```
In [7]: 1 # Declare an instance of the Softmax class.
2 # Weights are initialized to a random value.
3 # Note, to keep people's first solutions consistent, we are going to use a r
4
5 np.random.seed(1)
6
7 num_classes = len(np.unique(y_train))
8 num_features = X_train.shape[1]
9
10 softmax = Softmax(dims=[num_classes, num_features])
```

### Softmax loss

```
In [5]: 1 ## Implement the loss function of the softmax using a for loop over
2 # the number of examples
3
4 loss = softmax.loss(X_train, y_train)
```

```
In [6]: 1 print(loss)
```

```
2.3277607028048966
```

## Question:

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

## Answer:

Because there are a total of 10 classes, an untrained classifier will guess 0.1 probability for each class. Therefore, for each sample we would expect an error of  $-\log(0.1) = 2.3$ . Therefore, the mean error is also expected to be 2.3.

### Softmax gradient

In [7]:

```
1  ## Calculate the gradient of the softmax loss in the Softmax class.
2  # For convenience, we'll write one function that computes the loss
3  #   and gradient together, softmax.loss_and_grad(X, y)
4  # You may copy and paste your loss code from softmax.loss() here, and then
5  #   use the appropriate intermediate values to calculate the gradient.
6
7  loss, grad = softmax.loss_and_grad(X_dev,y_dev)
8
9  # Compare your gradient to a gradient check we wrote.
10 # You should see relative gradient errors on the order of 1e-07 or less if y
11 softmax.grad_check_sparse(X_dev, y_dev, grad)
```

```
numerical: -1.266499 analytic: -1.266499, relative error: 1.289325e-08
numerical: 0.056842 analytic: 0.056842, relative error: 4.306434e-07
numerical: 0.449315 analytic: 0.449315, relative error: 2.518900e-08
numerical: 0.935716 analytic: 0.935716, relative error: 1.540367e-08
numerical: -0.512110 analytic: -0.512110, relative error: 9.848258e-08
numerical: 0.342088 analytic: 0.342088, relative error: 2.513967e-08
numerical: -0.967954 analytic: -0.967954, relative error: 6.087353e-08
numerical: -1.891557 analytic: -1.891557, relative error: 1.427829e-08
numerical: -0.247811 analytic: -0.247811, relative error: 6.318535e-08
numerical: -2.658215 analytic: -2.658215, relative error: 2.500902e-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]:

```
1  import time
```

```

In [9]: 1  ## Implement softmax.fast_loss_and_grad which calculates the loss and gradie
        2  # WITHOUT using any for loops.
        3
        4  # Standard loss and gradient
        5  tic = time.time()
        6  loss, grad = softmax.loss_and_grad(X_dev, y_dev)
        7  toc = time.time()
        8  print('Normal loss / grad_norm: {} / {} computed in {}s'.format(loss, np.linalg.norm(grad), toc - tic))
        9
       10  tic = time.time()
       11  loss_vectorized, grad_vectorized = softmax.fast_loss_and_grad(X_dev, y_dev)
       12  toc = time.time()
       13  print('Vectorized loss / grad: {} / {} computed in {}s'.format(loss_vectorized, np.linalg.norm(grad_vectorized), toc - tic))
       14
       15  # The losses should match but your vectorized implementation should be much faster
       16  print('difference in loss / grad: {} / {} '.format(loss - loss_vectorized, np.linalg.norm(grad - grad_vectorized)))
       17
       18  # You should notice a speedup with the same output.

```

Normal loss / grad\_norm: 2.3230138756030048 / 281.3263564766588 computed in 0.07181096076965332s

Vectorized loss / grad: 2.323013874528945 / 281.3263564766588 computed in 0.009247779846191406s

difference in loss / grad: 1.0740599520886462e-09 / 1.9427077332200875e-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:

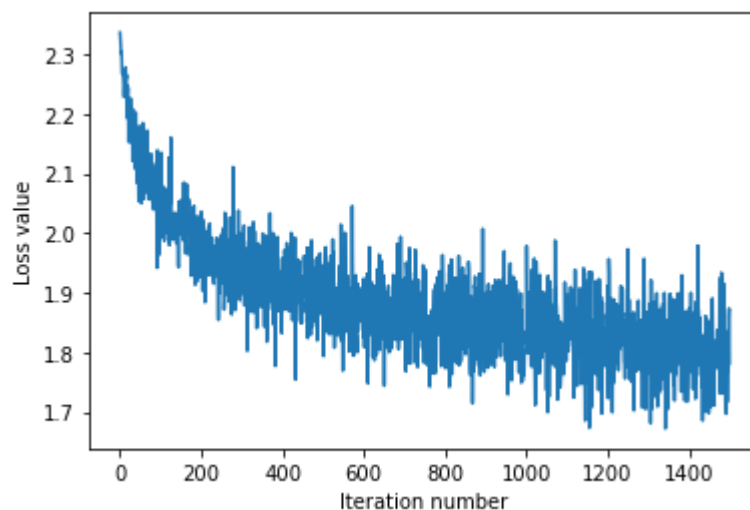
How should the softmax gradient descent training step differ from the svm training step, if at all?

### Answer:

Due to the change in loss function, the value of the gradients that are used to update the coefficient will change

```
In [46]: 1 # Implement softmax.train() by filling in the code to extract a batch of dat
2 # and perform the gradient step.
3 import time
4
5
6 tic = time.time()
7 loss_hist = softmax.train(X_train, y_train, learning_rate=1e-7,
8                           num_iters=1500, verbose=True)
9 toc = time.time()
10 print('That took {}s'.format(toc - tic))
11
12 plt.plot(loss_hist)
13 plt.xlabel('Iteration number')
14 plt.ylabel('Loss value')
15 plt.show()
```

```
iteration 0 / 1500: loss 2.3365926606637544
iteration 100 / 1500: loss 2.0557222613850827
iteration 200 / 1500: loss 2.0357745120662813
iteration 300 / 1500: loss 1.9813348165609888
iteration 400 / 1500: loss 1.9583142443981612
iteration 500 / 1500: loss 1.8622653073541355
iteration 600 / 1500: loss 1.8532611454359382
iteration 700 / 1500: loss 1.8353062223725827
iteration 800 / 1500: loss 1.829389246882764
iteration 900 / 1500: loss 1.8992158530357484
iteration 1000 / 1500: loss 1.97835035402523
iteration 1100 / 1500: loss 1.8470797913532633
iteration 1200 / 1500: loss 1.8411450268664082
iteration 1300 / 1500: loss 1.7910402495792102
iteration 1400 / 1500: loss 1.8705803029382257
That took 6.137814283370972s
```



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

```
In [47]: 1  ## Implement softmax.predict() and use it to compute the training and testin
          2
          3  y_train_pred = softmax.predict(X_train)
          4  print('training accuracy: {}'.format(np.mean(np.equal(y_train,y_train_pred),
          5  y_val_pred = softmax.predict(X_val)
          6  print('validation accuracy: {}'.format(np.mean(np.equal(y_val, y_val_pred)),
```

```
training accuracy: 0.3811428571428571
validation accuracy: 0.398
```

## Optimize the softmax classifier

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

```
In [10]: 1  np.finfo(float).eps
```

```
Out[10]: 2.220446049250313e-16
```

In [11]:

```

1  # ===== #
2  # YOUR CODE HERE:
3  #   Train the Softmax classifier with different Learning rates and
4  #   evaluate on the validation data.
5  #   Report:
6  #       - The best Learning rate of the ones you tested.
7  #       - The best validation accuracy corresponding to the best validation er
8  #
9  #   Select the SVM that achieved the best validation error and report
10 #   its error rate on the test set.
11 # ===== #
12 softmax = Softmax(dims=[num_classes, num_features])
13 rates = [10**i for i in range(-10,0)]
14
15 valAccuracies = []
16 valLosses = []
17 for rate in rates:
18     softmax.train(X_train, y_train, learning_rate=rate, num_iters=1500, verb
19     valLoss,a = softmax.fast_loss_and_grad(X_val, y_val)
20     valPreds = softmax.predict(X_val)
21     valAcc = np.mean(np.equal(y_val, valPreds))
22     valAccuracies.append(valAcc)
23     valLosses.append(valLoss)
24     print('Current rate:',rate,'validation accuracy: {}'.format(valAcc),'val
25     print("Best validation loss so far {}".format(min(valLosses)))
26 print("-"*100)
27 bestIndex = np.argmin(valLosses)
28 bestLR = rates[bestIndex]
29 bestValLoss = valLosses[bestIndex]
30 bestValAcc = valAccuracies[bestIndex]
31 print("The best learning rate is %f.Best validation loss is %f. Best validat
32
33
34 softmax.train(X_train, y_train, learning_rate=bestLR, num_iters=1500, verbos
35 yTestPred = softmax.predict(X_test)
36 testAcc = np.mean(np.equal(y_test, yTestPred))
37 print('Error of the softmax classifier with best learning rate %f on the tes
38 # ===== #
39 # END YOUR CODE HERE
40 # ===== #
41

```

Current rate: 1e-10 validation accuracy: 0.13 validation Loss: 2.332963609714193

Best validation loss so far 2.332963609714193

Current rate: 1e-09 validation accuracy: 0.179 validation Loss: 2.2418688560194973

Best validation loss so far 2.2418688560194973

Current rate: 1e-08 validation accuracy: 0.302 validation Loss: 2.020144565019459

Best validation loss so far 2.020144565019459

Current rate: 1e-07 validation accuracy: 0.379 validation Loss: 1.8287151179489731

Best validation loss so far 1.8287151179489731

Current rate: 1e-06 validation accuracy: 0.413 validation Loss: 1.7466767713611



042

Best validation loss so far 1.7466767713611042

Current rate: 1e-05 validation accuracy: 0.316 validation Loss: 2.5812750026016156

Best validation loss so far 1.7466767713611042

Current rate: 0.0001 validation accuracy: 0.256 validation Loss: 13.730043826690776

Best validation loss so far 1.7466767713611042

Current rate: 0.001 validation accuracy: 0.254 validation Loss: 16.890411756297777

Best validation loss so far 1.7466767713611042

Current rate: 0.01 validation accuracy: 0.145 validation Loss: 16.716676320502845

Best validation loss so far 1.7466767713611042

Current rate: 0.1 validation accuracy: 0.181 validation Loss: 17.1312330918501

Best validation loss so far 1.7466767713611042

-----  
-----

The best learning rate is 0.000001. Best validation loss is 1.746677. Best validation accuracy is 0.413000

Error of the softmax classifier with best learning rate 0.000001 on the test Set Error Rate is 0.622000