## Homework 2

# 1   K Nearest Neighbors

## 1.1   Problem Statement

(20 points) k-nearest neighbors. Complete the k-nearest neighbors Jupyter notebook. The goal of this workbook is to give you experience with the CIFAR-10 dataset, training and evaluating a simple classifier, and k-fold cross validation. In the Jupyter notebook, we'll be using the CIFAR-10 dataset. Acquire this dataset by running:

```
wget http://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz
tar -xzvf cifar-10-python.tar.gz
rm cifar-10-python.tar.gz
```

## 1.2 Jupyter Results

# knn

January 26, 2021

### 0.1 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.

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.

### 0.2 Import the appropriate libraries

```python
[131]: 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 files.
       # see http://stackoverflow.com/questions/1907993/
        ↪autoreload-of-modules-in-ipython
       %load_ext autoreload
       %autoreload 2
```

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

```python
[132]: # Set the path to the CIFAR-10 data
       cifar10_dir = 'cifar-10-batches-py' # You need to update this line
       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)
```

1

```
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,)
```

[133]:
```
# 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()
```

```python
[134]: # 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)
```

# 1 K-nearest neighbors

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

```python
[135]: # Import the KNN class

       from nndl import KNN
```

```python
[136]: # 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)
```

## 1.1 Questions

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

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

## 1.2 Answers

(1) knn.train() just saves all the training data and labels from the training set CIFAR

(2) The pros are the it is simple and fast (O(1)). On the other hand, it is memory intensive because all the training data must be stored and it scales with the amount of training examples.

## 1.3 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.

```
[137]:  # Implement the function compute_distances() in the KNN class.
        # Do not worry about the input 'norm' for now; use the default definition of␣
        ↪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: 26.693132877349854
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

### 1.3.1 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.

```
[138]:  # Implement the function compute_L2_distances_vectorized() in the KNN class.
        # In this function, you ought to achieve the same L2 distance but WITHOUT any␣
        ↪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: 0.14873099327087402
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.

### 1.3.2 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.

```python
[139]: # 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.
       # ================================================================ #
       y_pred = knn.predict_labels(dists_L2_vectorized,1)
       num_incorrect = np.count_nonzero((y_pred-y_test))
       error = num_incorrect/len(y_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%.

## 2 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.

### 2.0.1 Create training and validation folds

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

```python
[140]: # 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]
# ================================================================ #

X_train_folds = np.array_split(X_train, num_folds)
# print(X_train_folds[0].shape)
y_train_folds = np.array_split(y_train, num_folds)
# print(y_train_folds[0].shape)


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

### 2.0.2 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.

```
[141]: 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.
# ================================================================ #

# preallocations
num_k = len(ks)
xerror = np.zeros(num_folds)
av_error = np.zeros(num_k)

# iterating through all the k values
for k in np.arange(num_k):
    # iterating through all the folds
    for i in np.arange(num_folds):
        # choosing the fold for validation
        X_val_fold = X_train_folds[i]
```

```python
        y_val_fold = y_train_folds[i]

        # preallocations
        X_train_sfold = []
        y_train_sfold = []

        # assigning the rest of the folds to be used for training
        for l in np.arange(num_folds):
            if l != i:
                X_train_sfold.extend(X_train_folds[l])
                y_train_sfold.extend(y_train_folds[l])

        # converting list to array
        X_train_sfold = np.array(X_train_sfold)
        y_train_sfold = np.array(y_train_sfold)
        # print(np.shape(X_train_sfold))

        knn.train(X=X_train_sfold, y=y_train_sfold)
        dists_fold = knn.compute_L2_distances_vectorized(X=X_val_fold)

        y_pred = knn.predict_labels(dists_fold,ks[k])
        num_incorrect = np.count_nonzero((y_pred-y_val_fold))
        xerror[i] = num_incorrect/len(y_val_fold)

    av_error[k] = 1/num_folds*np.sum(xerror)

min_id = np.argmin(av_error)
plt.plot(ks,av_error)
plt.title('Nearest Neighbors, k Optimization')
plt.xlabel('k')
plt.ylabel('cross-validation error')
print('Optimum k value = %d, with error = %f' %(ks[min_id],av_error[min_id]))

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

print('Computation time: %.2f'%(time.time()-time_start))
```
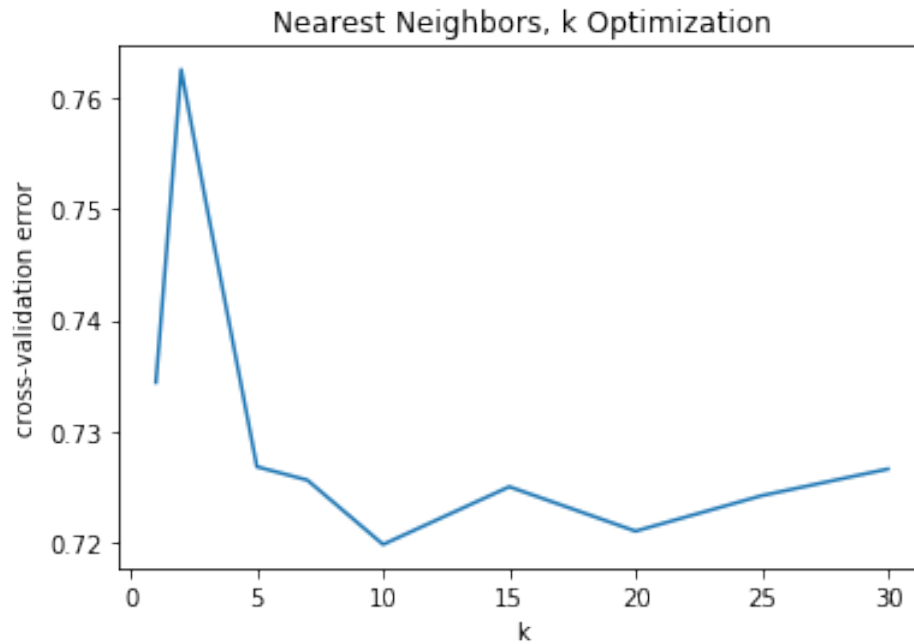
```
Optimum k value = 10, with error = 0.719800
Computation time: 19.98
```

7

Nearest Neighbors, k Optimization

## 2.1 Questions:

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

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

## 2.2 Answers:

(1) k=10 is the best amongst the tested k values.

(2) The cross-validation error for this value is 0.71980.

### 2.2.1 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.

```
[145]: 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]

# ================================================================ #
```

```python
# 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.
# ================================================================ #

# preallocations
num_norm = len(norms)
xerror = np.zeros(num_folds)
av_error = np.zeros(num_norm)
# k = 10
k_opt = ks[min_id]
# print(k_opt)

# iterating through all the k values
for k in np.arange(num_norm):
    # iterating through all the folds
    for i in np.arange(num_folds):
        #knn = KNN()

        # choosing the fold for validation
        X_val_fold = X_train_folds[i]
        y_val_fold = y_train_folds[i]

        # preallocations
        X_train_sfold = []
        y_train_sfold = []

        # assigning the rest of the folds to be used for training
        for l in np.arange(num_folds):
            if l != i:
                X_train_sfold.extend(X_train_folds[l])
                y_train_sfold.extend(y_train_folds[l])

        # converting list to array
        X_train_sfold = np.array(X_train_sfold)
        y_train_sfold = np.array(y_train_sfold)
        # print(np.shape(X_train_sfold))

        knn.train(X=X_train_sfold, y=y_train_sfold)
        dists_fold = knn.compute_distances(X = X_val_fold, norm=norms[k])
```

9

```
        y_pred = knn.predict_labels(dists_fold,k_opt)
        num_incorrect = np.count_nonzero((y_pred-y_val_fold))
        xerror[i] = num_incorrect/len(y_val_fold)

    av_error[k] = 1/num_folds*np.sum(xerror)

norm_name = ["L1","L2","Linf"]

plt.plot(norm_name,av_error)
plt.title('Norm Optimization')
plt.xlabel('norm used')
plt.ylabel('cross-validation error')




# ================================================================ #
# END YOUR CODE HERE
# ================================================================ #
print('Computation time: %.2f'%(time.time()-time_start))
```
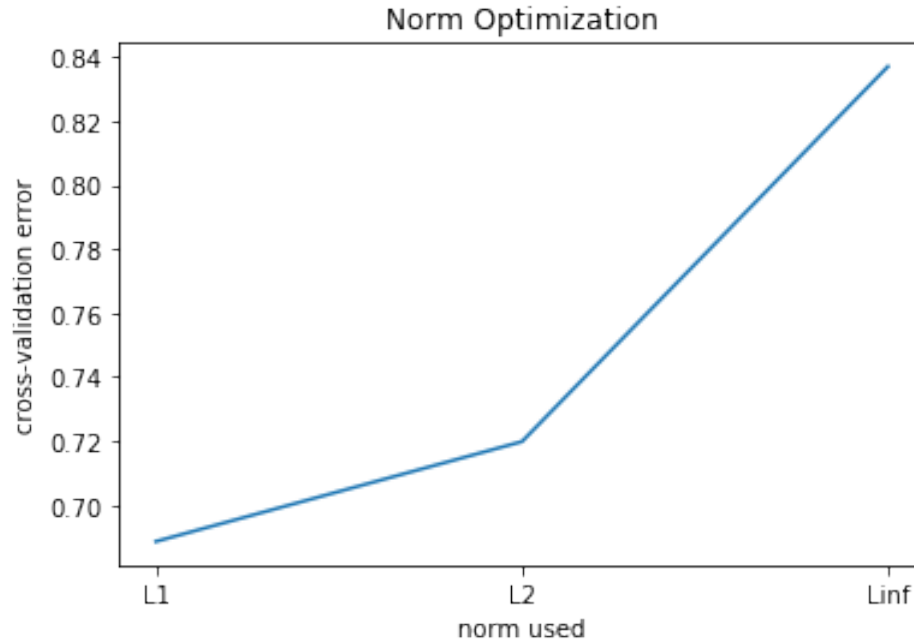
10
Computation time: 555.03



Norm Optimization

### 2.3 Questions:

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

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

### 2.4 Answers:

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

(2) The cross validation error for the L1 norm and k=10 is equal to 0.6886

## 3 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.

```
[148]:  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)

        norm_opt = norms[1]
        dists_L1= knn.compute_distances(X = X_test, norm = norm_opt)
        y_pred = knn.predict_labels(dists_L1,k_opt)
        num_incorrect = np.count_nonzero((y_pred-y_test))
        error = num_incorrect/len(y_test)

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

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

```
Error rate achieved: 0.718
```

### 3.1 Question:

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

### 3.2 Answer:

It improved from 0.726 to 0.718 which is about a 1% improvement

## 1.3 knn.py

```python
import numpy as np
import pdb

"""
This code was based off of code from cs231n at Stanford University, and ...
    modified for ECE C147/C247 at UCLA.
"""

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.
```

```
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          dists[i,j] = norm(X[i]-self.X_train[j])

52

53          # ...
                ================================================================ ...
```

```
                    #
54          # END YOUR CODE HERE
55          # ...
                    ================================================================  ...
                    #
56
57      return dists
58
59  def compute_L2_distances_vectorized(self, X):
60      """
61      Compute the distance between each test point in X and each training ...
              point
62      in self.X_train WITHOUT using any for loops.
63
64      Inputs:
65      - X: A numpy array of shape (num_test, D) containing test data.
66
67      Returns:
68      - dists: A numpy array of shape (num_test, num_train) where ...
               dists[i, j]
69        is the Euclidean distance between the ith test point and the jth ...
                 training
70        point.
71      """
72      num_test = X.shape[0]
73      num_train = self.X_train.shape[0]
74      dists = np.zeros((num_test, num_train))
75
76      # ================================================================ #
77      # YOUR CODE HERE:
78      #   Compute the L2 distance between the ith test point and the jth
79      #   training point and store the result in dists[i, j].  You may
80      #   NOT use a for loop (or list comprehension).  You may only use
```

15

```
81      #    numpy operations.
82      #
83      #    HINT: use broadcasting.  If you have a shape (N,1) array and
84      #    a shape (M,) array, adding them together produces a shape (N, M)
85      #    array.
86      # ================================================================ #
87
88      # euclidean norm = sqrt((xi-yj)^2) -> sqrt(xi^2+yj^2-2*xiyj)
89      # shape sizes
90      # (500,1)+(5000,) - 2* (500,5000)
91      # note axis 1 sums along rows
92      dists = np.sqrt(np.sum(X**2,axis = ...
            1,keepdims=True)+np.sum(self.X_train**2,axis=1)-2*np.dot(X,self.X_train.T))
93
94      # ================================================================ #
95      # END YOUR CODE HERE
96      # ================================================================ #
97
98      return dists
99
100
101  def predict_labels(self, dists, k=1):
102      """
103      Given a matrix of distances between test points and training points,
104      predict a label for each test point.
105
106      Inputs:
107      - dists: A numpy array of shape (num_test, num_train) where ...
            dists[i, j]
108        gives the distance betwen the ith test point and the jth training ...
            point.
109
110      Returns:
```

16

```
111        - y: A numpy array of shape (num_test,) containing predicted labels ...
              for the
112          test data, where y[i] is the predicted label for the test point X[i].
113        """
114        num_test = dists.shape[0]
115        y_pred = np.zeros(num_test)
116        for i in np.arange(num_test):
117          # A list of length k storing the labels of the k nearest ...
                neighbors to
118          # the ith test point.
119          closest_y = []
120          # ...
                 ================================================================= ...
                 #
121          # YOUR CODE HERE:
122          #   Use the distances to calculate and then store the labels of
123          #   the k-nearest neighbors to the ith test point.  The function
124          #   numpy.argsort may be useful.
125          #
126          #   After doing this, find the most common label of the k-nearest
127          #   neighbors.  Store the predicted label of the ith training example
128          #   as y_pred[i].  Break ties by choosing the smaller label.
129          # ...
                 ================================================================= ...
                 #
130
131          sortedIdxs = np.argsort(dists[i])
132          closest_y = self.y_train[sortedIdxs[:k]]
133          y_pred[i] = np.argmax(np.bincount(closest_y))
134
135          # ...
                 ================================================================= ...
                 #
```

```
136        # END YOUR CODE HERE
137        # ...
           =================================================================  ...
           #

138
139     return y_pred
```

# 2 Support Vector Machine

## 2.1 Problem Statement

(40 points) Support vector machine. Complete the SVM Jupyter notebook. Print out the entire workbook and related code sections in svm.py, then submit them as a pdf to gradescope.

## 2.2 Jupyter Results

svm

January 26, 2021

### 0.1 This is the svm workbook for ECE C147/C247 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.

### 0.2 Importing libraries and data setup

```python
[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 files.
     # see http://stackoverflow.com/questions/1907993/
      ↪autoreload-of-modules-in-ipython
     %load_ext autoreload
     %autoreload 2
```

```python
[511]: # Set the path to the CIFAR-10 data
       cifar10_dir = 'cifar-10-batches-py' # You need to update this line
       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)
```
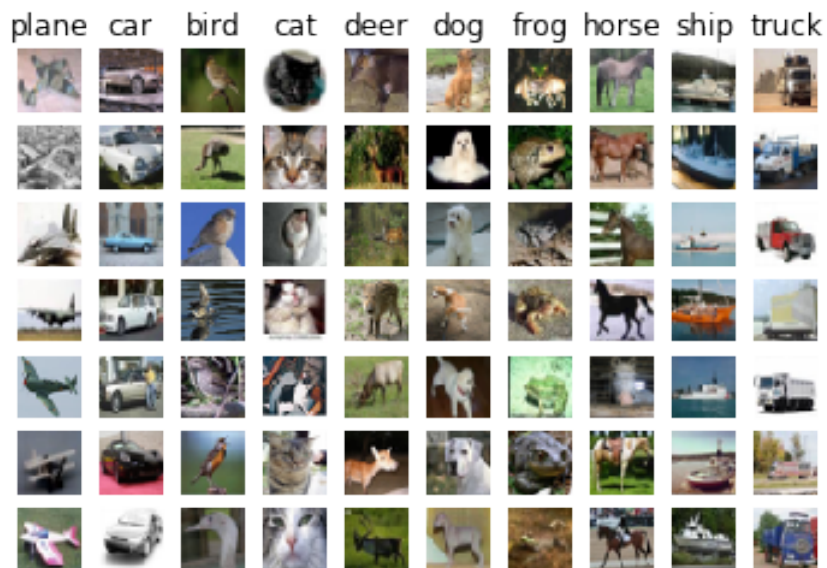
1

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

[512]:
```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()
```

```python
[513]:  # 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)
```

3

```
Dev labels shape:  (500,)
```

[514]:
```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)
```
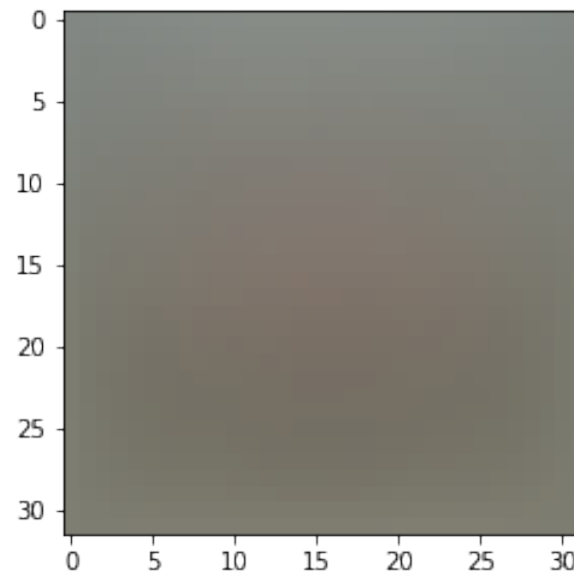
```
Training data shape:  (49000, 3072)
Validation data shape:  (1000, 3072)
Test data shape:  (1000, 3072)
dev data shape:  (500, 3072)
```

[515]:
```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 mean
 ↪image
plt.show()
```

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

```
[516]:  # 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
```

```
[517]:  # third: append the bias dimension of ones (i.e. bias trick) so that our SVM
        # 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)

## 0.3   Question:

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

## 0.4   Answer:

(1) For SVM mean-subtraction is performed because we want to "center" the data while in KNN we measure distances between points in order to distinguish classes so mean-subtraction would not do anything.

5

## 0.5   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.

```
[518]: from nndl.svm import SVM
```

```
[519]: # 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**

```
[520]: ## 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 15569.977915410187.

```
[ ]:
```

**SVM gradient**

```
[521]: ## 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: -3.681664 analytic: -3.681663, relative error: 2.542594e-08
numerical: 4.343144 analytic: 4.343143, relative error: 1.075819e-07
numerical: -3.260883 analytic: -3.260883, relative error: 9.346886e-09
```

6

```
numerical: 18.110456 analytic: 18.110456, relative error: 4.575248e-09
numerical: 10.624520 analytic: 10.624521, relative error: 2.793705e-08
numerical: 5.241199 analytic: 5.241198, relative error: 3.918448e-08
numerical: 5.384190 analytic: 5.384190, relative error: 5.281280e-08
numerical: -3.250758 analytic: -3.250758, relative error: 2.035146e-09
numerical: -6.618687 analytic: -6.618686, relative error: 3.693142e-08
numerical: -17.852199 analytic: -17.852198, relative error: 5.444536e-09
```

## 0.6  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.

```
[522]: import time
```

```
[523]: ## 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.linalg.
        ↪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_vectorized,␣
        ↪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, np.
        ↪linalg.norm(grad - grad_vectorized)))

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

```
Normal loss / grad_norm: 15816.57339070797 / 2205.1819144473743 computed in
0.10318613052368164s
Vectorized loss / grad: 15816.573390707988 / 2205.1819144473743 computed in
0.0050106048583984375s
difference in loss / grad: -1.8189894035458565e-11 / 2.6413939701336207e-12
```
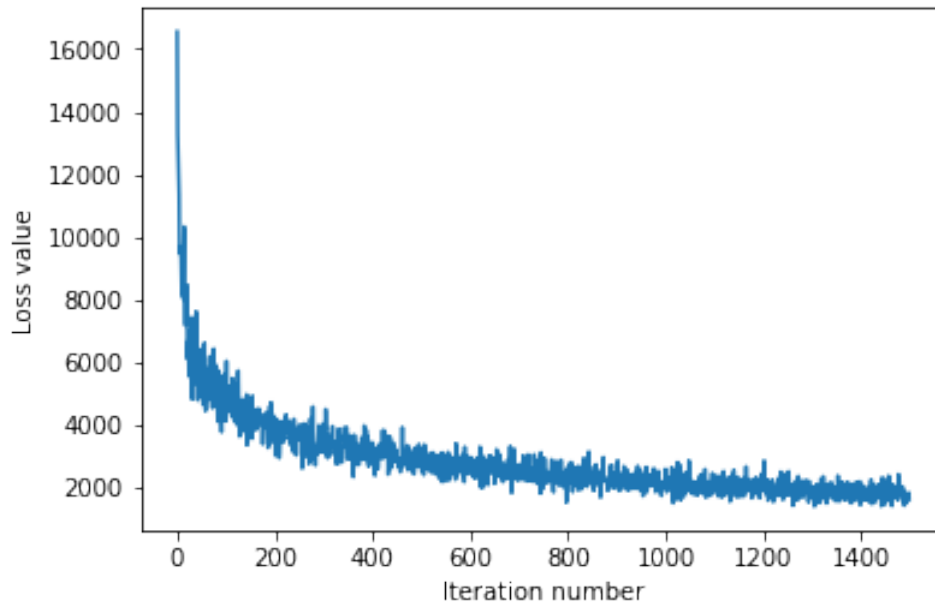
## 0.7  Stochastic gradient descent

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).

```
[524]: # 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 16557.38000190916
iteration 100 / 1500: loss 4701.089451272714
iteration 200 / 1500: loss 4017.333137942788
iteration 300 / 1500: loss 3681.9226471953625
iteration 400 / 1500: loss 2732.6164373988995
iteration 500 / 1500: loss 2786.6378424645045
iteration 600 / 1500: loss 2837.035784278268
iteration 700 / 1500: loss 2206.234868739933
iteration 800 / 1500: loss 2269.038824116981
iteration 900 / 1500: loss 2543.23781538592
iteration 1000 / 1500: loss 2566.6921357268257
iteration 1100 / 1500: loss 2182.0689059051633
iteration 1200 / 1500: loss 1861.1182244250447
iteration 1300 / 1500: loss 1982.9013858528256
iteration 1400 / 1500: loss 1927.5204158582117
That took 4.7113938331604s
```

### 0.7.1 Evaluate the performance of the trained SVM on the validation data.

```
[525]:  ## 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: 0.28530612244897957
validation accuracy: 0.3
```

## 0.8 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).

```
[526]:  # ================================================================ #
        # 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.
        #
```

9

```python
#   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
# ================================================================ #
epsilon = 5*np.logspace(-10,-1,num = 10)
val_acc = np.zeros_like(epsilon)

ep_opt = 0
val_opt = 0

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

    svm.train(X_train, y_train, learning_rate=epsilon[i],num_iters=1500,
 →verbose=False)
    y_val_pred = svm.predict(X_val)
    val_acc[i] = np.mean(np.equal(y_val, y_val_pred))

    if val_acc[i]> val_opt:
        val_opt = val_acc[i]
        ep_opt = epsilon[i]


print('Optimal epsilon = %4.3e with validation accuracy %4.3f'
 →%(ep_opt,val_opt))
plt.semilogx(epsilon,val_acc)
plt.xlabel('Learning Rate')
plt.ylabel('Validation Accuracy')
plt.show()

svm.train(X_train, y_train, learning_rate=ep_opt,num_iters=1500, verbose=False)
y_test_pred = svm.predict(X_test)
test_err = 1-np.mean(np.equal(y_test, y_test_pred))
print('Optimal epsilon = %4.3e with test error %4.3f' %(ep_opt,test_err))
# ================================================================ #
# END YOUR CODE HERE
# ================================================================ #
```
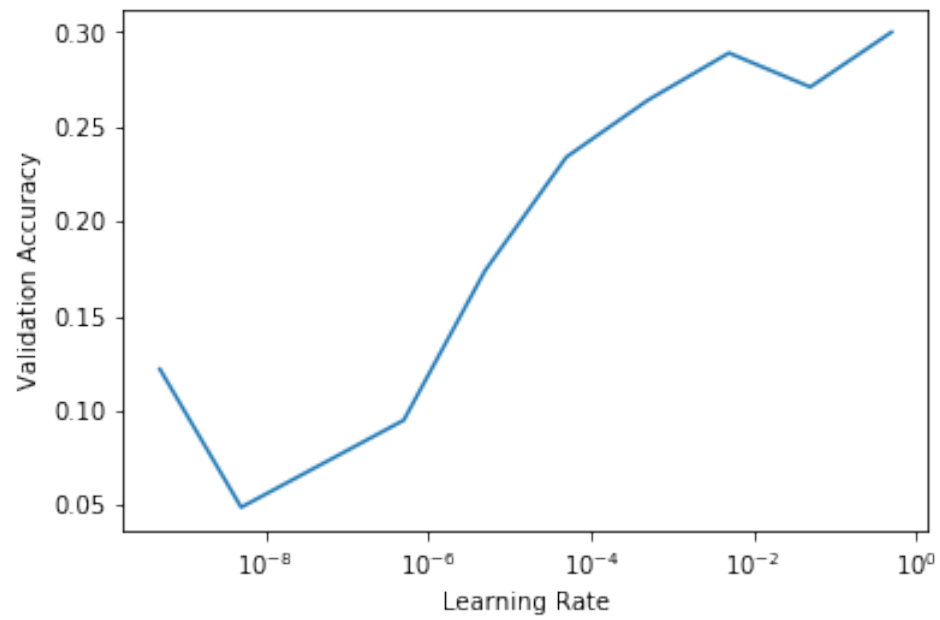
Optimal epsilon = 5.000e-01 with validation accuracy 0.300

Optimal epsilon = 5.000e-01 with test error 0.741

[ ]:

## 2.3   svm.py

```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 ECE C147/C247 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.
```

```python
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      # YOUR CODE HERE:
43        #   Calculate the normalized SVM loss, and store it as 'loss'.
44      #   (That is, calculate the sum of the losses of all the training
45      #   set margins, and then normalize the loss by the number of
46        #   training examples.)
47      # ================================================================ #
48        a = self.W.dot(X[i].T)  #(C,1)
49
50        for j in np.arange(num_classes):
51          if j != y[i]:
52            hinge = np.maximum(0,1+a[j]-a[y[i]])
53            loss += hinge
54      loss = 1/num_train*(loss)
55
56      # ================================================================ #
57      # END YOUR CODE HERE
58      # ================================================================ #
59
```

```
60      return loss

61

62   def loss_and_grad(self, X, y):
63      """
64      Same as self.loss(X, y), except that it also returns the gradient.

65

66      Output: grad -- a matrix of the same dimensions as W containing
67          the gradient of the loss with respect to W.
68      """

69

70      # compute the loss and the gradient
71      num_classes = self.W.shape[0]
72      num_train = X.shape[0]
73      loss = 0.0
74      grad = np.zeros_like(self.W)

75

76      for i in np.arange(num_train):
77      # ================================================================ #
78      # YOUR CODE HERE:
79        #   Calculate the SVM loss and the gradient.  Store the gradient in
80      #   the variable grad.
81      # ================================================================ #
82        a = self.W.dot(X[i].T) #(C,1)

83

84        for j in np.arange(num_classes):
85          zj = 1+a[j]-a[y[i]]

86

87          #defining indicator function
88          if zj > 0:
89            ind = 1
90          else:
91            ind = 0

92
```

```python
93              if j != y[i]:
94                  hinge = np.maximum(0,zj)
95                  loss += hinge
96                  grad[j] += ind*X[i]  #(1,D)
97                  grad[y[i]] -= ind*X[i]
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 ouptuts 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      # preallocations
142      a = X.dot(self.W.T) #size (N,C)
143      num_train = a.shape[0]
144      # num_class = a.shape[1]

145
146      a_ind = [np.arange(num_train),y]
147      hinge = np.maximum(np.zeros_like(a),np.ones_like(a)+(a.T-a[a_ind]).T)

148
149      # for j = yi hingeloss = 0
150      hinge[a_ind] = 0
151      loss = np.sum(hinge)
152      loss = 1/num_train*(loss)

153
154      # ================================================================= #
155      # END YOUR CODE HERE
156      # ================================================================= #
```

```
157

158

159

160         # ================================================================= #
161         # YOUR CODE HERE:
162         #   Calculate the SVM grad WITHOUT any for loops.
163         # ================================================================= #
164         # defining indicator function
165         ind = hinge #(N,C)
166         ind[ind>0] = 1

167

168         # summing through all the training samples
169         sumyi = np.sum(ind, axis =1)
170         # changing j = yi indices to -sum of j!=i
171         ind[a_ind] = -sumyi.T
172         grad = ind.T.dot(X) # (C,D)
173         grad = 1/num_train*(grad)

174

175         # ================================================================= #
176         # END YOUR CODE HERE
177         # ================================================================= #

178

179         return loss, grad

180

181    def train(self, X, y, learning_rate=1e-3, num_iters=100,
182              batch_size=200, verbose=False):
183        """
184        Train this linear classifier using stochastic gradient descent.

185

186        Inputs:
187        - X: A numpy array of shape (N, D) containing training data; there ...
              are N
188          training samples each of dimension D.
```

```
189    - y: A numpy array of shape (N,) containing training labels; y[i] = c
190      means that X[i] has label 0 ≤ c < C for C classes.
191    - learning_rate: (float) learning rate for optimization.
192    - num_iters: (integer) number of steps to take when optimizing
193    - batch_size: (integer) number of training examples to use at each ...
          step.
194    - verbose: (boolean) If true, print progress during optimization.
195
196    Outputs:
197    A list containing the value of the loss function at each training ...
          iteration.
198    """
199    num_train, dim = X.shape
200    num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K ...
          is number of classes
201
202    self.init_weights(dims=[np.max(y) + 1, X.shape[1]]) # initializes ...
          the weights of self.W
203
204    # Run stochastic gradient descent to optimize W
205    loss_history = []
206
207    for it in np.arange(num_iters):
208      X_batch = None
209      y_batch = None
210
211      # ...
              ===================================================================== ...
              #
212      # YOUR CODE HERE:
213      #   Sample batch_size elements from the training data for use in
214      #   gradient descent.  After sampling,
215      #     - X_batch should have shape: (dim, batch_size)
```

```
216        #      - y_batch should have shape: (batch_size,)
217        #   The indices should be randomly generated to reduce correlations
218        #   in the dataset.  Use np.random.choice.  It's okay to sample ...
              with
219        #   replacement.
220     # ...
           ================================================================ ...
             #
221     index = np.random.choice(num_train,batch_size)
222     X_batch = X[index]
223     y_batch = y[index]
224     # ...
           ================================================================ ...
             #
225     # END YOUR CODE HERE
226     # ...
           ================================================================ ...
             #
227
228     # evaluate loss and gradient
229     loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
230     loss_history.append(loss)
231
232     # ...
           ================================================================ ...
             #
233     # YOUR CODE HERE:
234     #   Update the parameters, self.W, with a gradient step
235     # ...
           ================================================================ ...
             #
236     self.W = self.W - learning_rate*grad
```

```
237            # ...
                 ===================================================================  ...
                 #
238        # END YOUR CODE HERE
239        # ...
                 ================================================================  ...
                 #
240
241        if verbose and it % 100 == 0:
242            print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
243
244        return loss_history
245
246    def predict(self, X):
247        """
248        Inputs:
249        - X: N x D array of training data. Each row is a D-dimensional point.
250
251        Returns:
252        - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
253            array of length N, and each element is an integer giving the ...
                   predicted
254            class.
255        """
256        y_pred = np.zeros(X.shape[1])
257
258
259        # ================================================================ #
260        # YOUR CODE HERE:
261        #    Predict the labels given the training data with the parameter ...
                 self.W.
262        # ================================================================ #
263        y_pred = np.argmax(X.dot(self.W.T),axis=1)
```

```
264        # ================================================================= #

265        # END YOUR CODE HERE

266        # ================================================================= #

267

268        return y_pred
```

# 3 Softmax

## 3.1 Problem Statement

(40 points) Softmax classifier. Complete the Softmax Jupyter notebook. Print out the entire workbook and related code sections in softmax.py, then submit them as a pdf to gradescope.

## 3.2 Jupyter Results

### softmax

#### January 26, 2021

### 0.1 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.

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.

```
[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
```

```
[48]: def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000,␣
      ↪num_dev=500):
          """
          Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
          it for the linear classifier. These are the same steps as we used for the
          SVM, but condensed to a single function.
          """
          # Load the raw CIFAR-10 data
          cifar10_dir = 'cifar-10-batches-py' # You need to update this line
          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]
```

1

```python
    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)
```

2

```
dev labels shape:  (500,)
```

## 0.2   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.

```
[49]: from nndl import Softmax
```

```
[50]: # 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**

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

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

```
[52]: print(loss)
```

```
2.3277607028048966
```

## 0.3   Question:

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

## 0.4   Answer:

-log(softmax(x)) = 2.3 which means softmax(x) ~ 0.10. This makes sense because there are ten classes meaning equal probability it could be one of the classes.

**Softmax gradient**

```
[53]: ## 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.
```

3

```
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: 0.037523 analytic: 0.037523, relative error: 3.387207e-07
numerical: 1.437461 analytic: 1.437461, relative error: 2.819630e-08
numerical: -0.726794 analytic: -0.726794, relative error: 9.528998e-09
numerical: 1.543500 analytic: 1.543500, relative error: 1.314298e-08
numerical: -0.511871 analytic: -0.511871, relative error: 1.320853e-07
numerical: 1.243152 analytic: 1.243152, relative error: 1.538739e-08
numerical: -0.451246 analytic: -0.451246, relative error: 1.064016e-07
numerical: -1.734786 analytic: -1.734786, relative error: 9.860642e-09
numerical: 0.414030 analytic: 0.414030, relative error: 6.750906e-09
numerical: -1.058630 analytic: -1.058630, relative error: 2.340584e-08
```

## 0.5 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.

```
[54]: import time
```

```
[55]: ## Implement softmax.fast_loss_and_grad which calculates the loss and gradient
      #     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.linalg.
       ↪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_vectorized,
       ↪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, np.
       ↪linalg.norm(grad - grad_vectorized)))

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

4

```
Normal loss / grad_norm: 2.3613109117810556 / 321.6294577431447 computed in
0.08613395690917969s
Vectorized loss / grad: 2.3613109117810542 / 321.62945774314477 computed in
0.004011392593383789s
difference in loss / grad: 1.3322676295501878e-15 /2.272637193003094e-13
```

## 0.6  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).

## 0.7  Question:

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

## 0.8  Answer:

The training step does not differ between the two methods besides the fact they have different loss functions
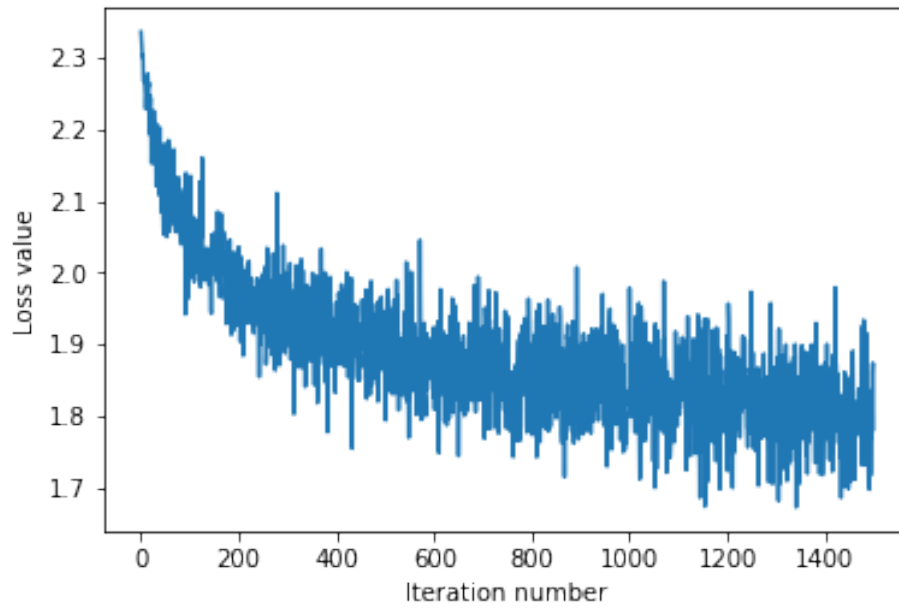
```
[56]: # Implement softmax.train() by filling in the code to extract a batch of data
      # and perform the gradient step.
      import time


      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 2.3365926606637544
iteration 100 / 1500: loss 2.0557222613850827
iteration 200 / 1500: loss 2.035774512066282
iteration 300 / 1500: loss 1.9813348165609888
iteration 400 / 1500: loss 1.9583142443981612
iteration 500 / 1500: loss 1.8622653073541355
iteration 600 / 1500: loss 1.8532611454359385
iteration 700 / 1500: loss 1.835306222372583
iteration 800 / 1500: loss 1.829389246882764
iteration 900 / 1500: loss 1.8992158530357484
iteration 1000 / 1500: loss 1.97835035402523
iteration 1100 / 1500: loss 1.8470797913532633
```

5

```
iteration 1200 / 1500: loss 1.8411450268664085
iteration 1300 / 1500: loss 1.7910402495792102
iteration 1400 / 1500: loss 1.870580302938226
That took 4.511507749557495s
```



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

```
[57]: ## 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: 0.3811428571428571
validation accuracy: 0.398
```

### 0.9 Optimize the softmax classifier

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

```
[58]: np.finfo(float).eps
```

6

[58]: `2.220446049250313e-16`

[59]:
```python
# ================================================================ #
# 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 error.
#
#   Select the SVM that achieved the best validation error and report
#     its error rate on the test set.
# ================================================================ #
epsilon = np.logspace(-10,-1, num=10)
val_acc = np.zeros_like(epsilon)

ep_opt = 0
val_opt = 0

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

    softmax.train(X_train, y_train, learning_rate=epsilon[i],num_iters=1500,␣
 ↪verbose=False)
    y_val_pred = softmax.predict(X_val)
    val_acc[i] = np.mean(np.equal(y_val, y_val_pred))

    if val_acc[i]> val_opt:
        val_opt = val_acc[i]
        ep_opt = epsilon[i]


print('Optimal epsilon = %4.3e with validation accuracy %4.3f'␣
 ↪%(ep_opt,val_opt))
plt.semilogx(epsilon,val_acc)
plt.xlabel('Learning Rate')
plt.ylabel('Validation Accuracy')
plt.show()

softmax.train(X_train, y_train, learning_rate=ep_opt,num_iters=1500,␣
 ↪verbose=False)
y_test_pred = softmax.predict(X_test)
test_err = 1- np.mean(np.equal(y_test, y_test_pred))
print('Optimal epsilon = %4.3e with test error %4.3f' %(ep_opt,test_err))


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

```
# =================================================================== #
```

Optimal epsilon = 1.000e-06 with validation accuracy 0.408



Optimal epsilon = 1.000e-06 with test error 0.598

[ ]:

## 3.3   softmax.py

```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        a = self.W.dot(X.T) #size (C,N)
44        num_sample = a.shape[1]
45        sumL = 0
46
47        # iterating through all the training samples (columns of a)
48        for i in np.arange(num_sample):
49          # to avoid overflow, normalize softmax func by logk = -maxai(x) ...
                 (maximum value in the column)
50          logk= -np.amax(a[:,i])
51
52          # y[i] is the class that the training image belongs to
53          # want to extract the score for the first x(i) for class y[i]
54          smax = np.exp(a[y[i],i]+logk)/np.sum(np.exp(a[:,i]+logk))
55          L_i = -np.log(smax)
56
57          sumL = sumL + L_i
58        loss = 1/num_sample*(sumL)
59
```

```python
60      # =================================================================== #
61      # END YOUR CODE HERE
62      # =================================================================== #
63
64      return loss
65
66  def loss_and_grad(self, X, y):
67      """
68      Same as self.loss(X, y), except that it also returns the gradient.
69
70      Output: grad -- a matrix of the same dimensions as W containing
71          the gradient of the loss with respect to W.
72      """
73
74      # Initialize the loss and gradient to zero.
75      loss = 0.0
76      grad = np.zeros_like(self.W)
77
78      # =================================================================== #
79      # YOUR CODE HERE:
80      #   Calculate the softmax loss and the gradient. Store the gradient
81      #   as the variable grad.
82      # =================================================================== #
83
84      # calculating loss
85      a = self.W.dot(X.T) #size (C,N)
86      num_sample = a.shape[1]
87      num_class = a.shape[0]
88
89      # iterating through all the training samples (columns of a)
90      for i in np.arange(num_sample):
91          # avoiding overflow
92          logk= -np.amax(a[:,i])
```

```
93        a[:,i] += logk

94

95        # iterating through classes

96        for j in np.arange(num_class):

97          smax = np.exp(a[j,i])/np.sum(np.exp(a[:,i]))

98          if y[i] == j:

99            # want to take grad of g = -log(smax)

100           # dg/dsmax = -1/smax

101           # dsmax/da = smax(1-smax)

102           # da/dw = x

103           # grad[j,:] += -1/smax*smax*(1-smax)*X[i,:]

104           grad[j,:] += (smax-1)*X[i,:]

105

106         else:

107           # dg/dsmax = -1/smax

108           # dsmax/da = -smax(smax)

109           # da/dw = x

110           # grad[j,:] += -1/smax*(-smax*smax)*X[i,:]

111           grad[j,:] += smax*X[i,:]

112

113     loss = self.loss(X,y)

114

115     # calculating gradient

116     # (C,D)

117     grad = grad/num_sample

118

119     # ================================================================ #

120     # END YOUR CODE HERE

121     # ================================================================ #

122

123     return loss, grad

124

125   def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
```

```
126     """
127     sample a few random elements and only return numerical
128     in these dimensions.
129     """
130
131     for i in np.arange(num_checks):
132       ix = tuple([np.random.randint(m) for m in self.W.shape])
133
134       oldval = self.W[ix]
135       self.W[ix] = oldval + h # increment by h
136       fxph = self.loss(X, y)
137       self.W[ix] = oldval - h # decrement by h
138       fxmh = self.loss(X,y) # evaluate f(x - h)
139       self.W[ix] = oldval # reset
140
141       grad_numerical = (fxph - fxmh) / (2 * h)
142       grad_analytic = your_grad[ix]
143       rel_error = abs(grad_numerical - grad_analytic) / ...
              (abs(grad_numerical) + abs(grad_analytic))
144       print('numerical: %f analytic: %f, relative error: %e' % ...
              (grad_numerical, grad_analytic, rel_error))
145
146   def fast_loss_and_grad(self, X, y):
147     """
148     A vectorized implementation of loss_and_grad. It shares the same
149     inputs and ouptuts as loss_and_grad.
150     """
151     loss = 0.0
152     grad = np.zeros(self.W.shape) # initialize the gradient as zero
153
154     # ================================================================ #
155     # YOUR CODE HERE:
156       #   Calculate the softmax loss and gradient WITHOUT any for loops.
```

```python
157      # ================================================================ #
158
159      a = X.dot(self.W.T) #size (N,C)
160      num_sample = a.shape[0]
161
162      #  normalize by max in each class
163      logk= -np.amax(a, axis =1, keepdims = True)
164      # creating logk vector into matrix of repeated column value
165      a = a + np.tile(logk,(1,a.shape[1]))
166
167      # defining softmax function
168      smax = np.exp(a)/np.sum(np.exp(a), axis =1, keepdims = True)
169      # y is a list of classes that x belongs to [1 c 3 ... c], want to ...
             select all the samples of xi that correspond to yi
170      a_ind = [np.arange(num_sample),y]
171      loss = np.sum(-np.log(smax[a_ind]))
172      loss = 1/num_sample*loss
173
174      # creating indicator function
175      ind = np.zeros_like(smax)
176      ind[a_ind] = 1
177
178      # recall for smax yi equal to j dg/df*df/dz = (smax-1) otherwise ...
             just smax
179      dgdz = smax - ind
180      grad = dgdz.T.dot(X)
181      grad = 1/num_sample*grad
182      # ================================================================ #
183      # END YOUR CODE HERE
184      # ================================================================ #
185
186      return loss, grad
187
```

```
188    def train(self, X, y, learning_rate=1e-3, num_iters=100,
189              batch_size=200, verbose=False):
190      """
191      Train this linear classifier using stochastic gradient descent.
192
193      Inputs:
194      - X: A numpy array of shape (N, D) containing training data; there ...
              are N
195        training samples each of dimension D.
196      - y: A numpy array of shape (N,) containing training labels; y[i] = c
197        means that X[i] has label 0 <= c < C for C classes.
198      - learning_rate: (float) learning rate for optimization.
199      - num_iters: (integer) number of steps to take when optimizing
200      - batch_size: (integer) number of training examples to use at each ...
              step.
201      - verbose: (boolean) If true, print progress during optimization.
202
203      Outputs:
204      A list containing the value of the loss function at each training ...
              iteration.
205      """
206      num_train, dim = X.shape
207      num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K ...
              is number of classes
208
209      self.init_weights(dims=[np.max(y) + 1, X.shape[1]]) # initializes ...
              the weights of self.W
210
211      # Run stochastic gradient descent to optimize W
212      loss_history = []
213
214      for it in np.arange(num_iters):
215        X_batch = None
```

```
216        y_batch = None

217

218        # ...
               ======================================================================  ...
               #
219        # YOUR CODE HERE:
220        #   Sample batch_size elements from the training data for use in
221        #   gradient descent.  After sampling,
222        #     - X_batch should have shape: (dim, batch_size)
223        #     - y_batch should have shape: (batch_size,)
224        #   The indices should be randomly generated to reduce correlations
225        #   in the dataset.  Use np.random.choice.  It's okay to sample with
226        #   replacement.
227        # ...
               ======================================================================  ...
               #
228        N = X.shape[0]
229        index = np.random.choice(N,batch_size)
230        X_batch = X[index]
231        y_batch = y[index]
232        # ...
               ======================================================================  ...
               #
233        # END YOUR CODE HERE
234        # ...
               ======================================================================  ...
               #

235

236        # evaluate loss and gradient
237        loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
238        loss_history.append(loss)

239
```

```
240        # ...
              ================================================================ ...
              #
241        # YOUR CODE HERE:
242        #   Update the parameters, self.W, with a gradient step
243        # ...
              ================================================================ ...
              #
244        self.W = self.W - learning_rate*grad
245
246          # ...
                 ================================================================ ...
                 #
247        # END YOUR CODE HERE
248        # ...
              ================================================================ ...
              #
249
250        if verbose and it % 100 == 0:
251          print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
252
253      return loss_history
254
255    def predict(self, X):
256      """
257      Inputs:
258      - X: N x D array of training data. Each row is a D-dimensional point.
259
260      Returns:
261      - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
262          array of length N, and each element is an integer giving the ...
                 predicted
263          class.
```

```
264        """

265        y_pred = np.zeros(X.shape[1])

266        # ================================================================ #

267        # YOUR CODE HERE:

268        #    Predict the labels given the training data.

269        # ================================================================ #

270        # find maximum score out of all the classes and return index

271        y_pred = np.argmax(X.dot(self.W.T),axis=1)

272        # ================================================================ #

273        # END YOUR CODE HERE

274        # ================================================================ #

275

276        return y_pred
```