

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

## Assignment #2

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

Please print out the workbook entirely when completed.

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

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

## Import the appropriate libraries

```
In [318]: 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:

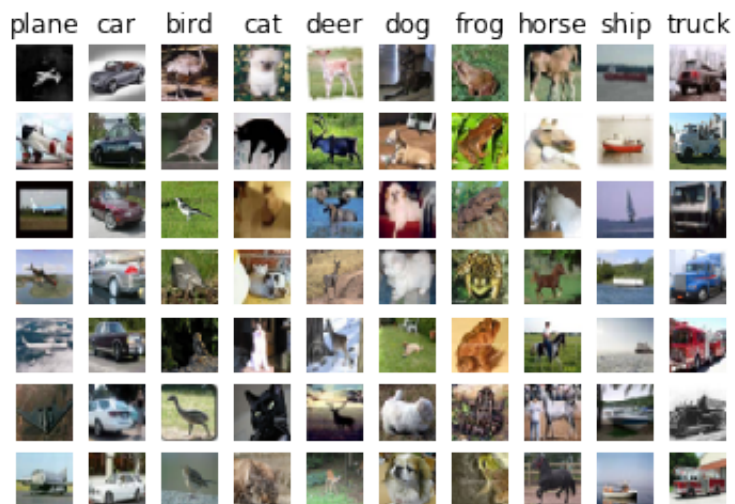
```
%reload_ext autoreload
```

```
In [319]: # Set the path to the CIFAR-10 data
cifar10_dir = 'cifar-10-batches-py'
X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

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

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

```
In [320]: # 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', 'horse', 'ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y) #return indices of the non-zero elements of the input array
    idxs = np.random.choice(idxs, samples_per_class, replace=False) #generates a random sample from a given 1D array
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()
```



```
In [321]: # Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = list(range(num_training))
X_train = X_train[mask]
y_train = y_train[mask]

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

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

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

## K-nearest neighbors

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

```
In [322]: # Import the KNN class

from nn1 import KNN
```

```
In [323]: # 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)

# print(np.array_equal(knn.compute_L2_distances_vectorized(X_test), knn.compute_distances(X_test)))
#knn.compute_L2_distances_vectorized(X_test)

#knn.predict_labels(dists=knn.compute_L2_distances_vectorized(X_test))
```

## 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) Function `knn.train()` takes two objects as inputs: a matrix of `X_train` as the training data `self.X_train`, and the corresponding labels `y_train` as the `self.y_train` to be called in the functions.

(2) Pros: no training time and thus simple. Cons: requires caching the entire training set, which could be impractical if large, is computationally expensive on testing new data, the curse of dimensionality may be at play and the data representation is very important.

## 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 [324]: # Implement the function compute_distances() in the KNN class.  
# Do not worry about the input 'norm' for now; use the default definit  
ion 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: 57.70186114311218  
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 [325]: # Implement the function compute_L2_distances_vectorized() in the KNN  
class.  
# In this function, you ought to achieve the same L2 distance but WITH  
OUT 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 (sh  
ould be 0): {}'.format(np.linalg.norm(dists_L2 - dists_L2_vectorized,  
'fro')))  
  
Time to run code: 0.4430050849914551  
Difference in L2 distances between your KNN implementations (should  
be 0): 0.0
```

## Speedup

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

## Implementing the prediction

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

```
In [326]: # 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_test_pred = knn.predict_labels(dists=knn.compute_L2_distances_vectorized(X_test))
#print(y_test_pred)
#print(y_test)
error -= float(np.sum(y_test_pred == y_test))/num_test

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

print(error)

0.726
```

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

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

## Optimizing KNN hyperparameters

In this section, we'll take the KNN classifier that you have constructed and perform cross-validation to choose a best value of  $k$ , as well as a best choice of norm. In k-fold cross-validation, the original sample is randomly partitioned into  $k$  equal size subsamples. Of the  $k$  subsamples, a single subsample is retained as the validation data for testing the model, and the remaining  $k-1$  subsamples are used as training data. The cross-validation process is then repeated  $k$  times (the folds), with each of the  $k$  subsamples used exactly once as the validation data. The  $k$  results from the folds can then be averaged (or otherwise combined) to produce a single estimation. The advantage of this method is that all observations are used for both training and validation, and each observation is used for validation exactly once.

### Create training and validation folds

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

```

In [236]: # Create the dataset folds for cross-validation.
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]
# ===== #
cv_idx = np.arange(num_training)
np.random.shuffle(cv_idx)
ind = np.array_split(cv_idx, num_folds)
for i in ind:
    X_train_folds.append(X_train[i])
    y_train_folds.append(y_train[i])
    #print(X_train_folds)
    #print(y_train_folds)

#X_train_folds_noshuffle = np.array_split(X_train, num_folds)
# ===== #
# END YOUR CODE HERE
# ===== #

```

## 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 [237]: 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.
# ===== #
k_error = {}
e = []
for k in ks:
    k_error[k] = []
    for f in range(num_folds):
        X_test = X_train_folds[f]
        y_test = y_train_folds[f]
        X_train = np.concatenate(X_train_folds[:f] + X_train_folds[(f+
1):])
        y_train = np.concatenate(y_train_folds[:f] + y_train_folds[(f+
1):])

        knn.train(X_train, y_train)

        dists = knn.compute_L2_distances_vectorized(X_test)
        y_test_pred = knn.predict_labels(dists, k=k)

        num_correct = np.sum(y_test_pred == y_test)
        error = 1 - float(num_correct) / X_test.shape[0]
        k_error[k].append(error)
#for k in k_error:
#    for avg_error in k_error[k]:
#        print (k,avg_error)
    avg_error = float(np.sum(k_error[k])/num_folds)
    print ('k =', k, 'Average error=', avg_error)
    e.append(avg_error)

plt.plot(ks, e, label = 'k-fold cross validation error')
plt.xlabel('k')
plt.ylabel('Cross-validation Error')
plt.grid()
plt.show()

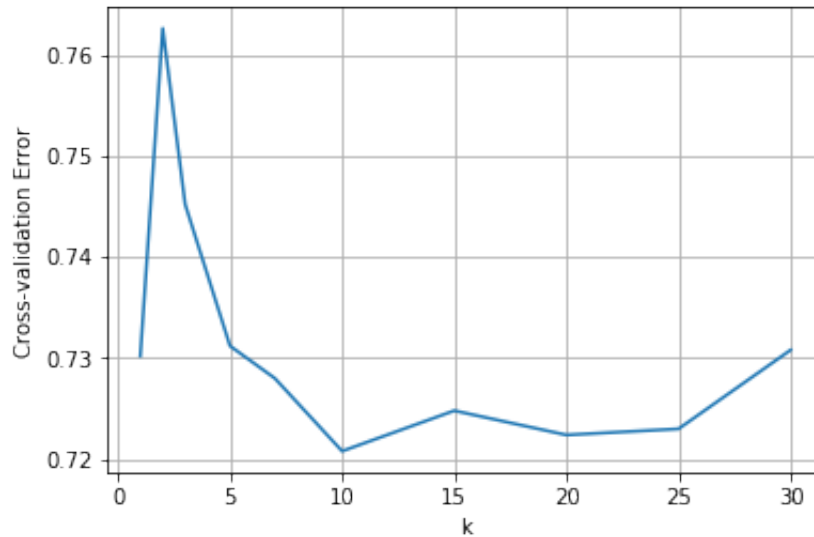
# ===== #
# END YOUR CODE HERE
# ===== #

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

```



k = 1 Average error= 0.7302000000000001  
k = 2 Average error= 0.7626  
k = 3 Average error= 0.7452  
k = 5 Average error= 0.7312  
k = 7 Average error= 0.728  
k = 10 Average error= 0.7208  
k = 15 Average error= 0.7248  
k = 20 Average error= 0.7224  
k = 25 Average error= 0.7230000000000001  
k = 30 Average error= 0.7308



Computation time: 47.39

## 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)  $k = 10$  is the best among the tested  $k$ 's, with the lowest cross-validation error rate if I do not shuffle the training data prior to splitting. If shuffled/used randomization, according to the results above, the best  $k$  can vary. In the example above, still  $k = 10$ .
- (2) 0.7198 for  $k=10$  without shuffling data (`#np.random.shuffle`), 0.7208 for  $k = 10$  in the run above based on shuffled training data.

## 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 [304]: time_start = time.time()

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

# ===== #
# YOUR CODE HERE:
# Calculate the cross-validation error for each norm in norms, testing
# the trained model on each of the 5 folds. Average these errors
# together and make a plot of the norm used vs the cross-validation
# error
# Use the best cross-validation k from the previous part.
#
# Feel free to use the compute_distances function. We're testing just
# three norms, but be advised that this could still take some time.
# You're welcome to write a vectorized form of the L1- and Linf- norms
# to speed this up, but it is not necessary.
# ===== #
Error_avg = []
for norm in norms:
    errorsum = 0
    t = time.time()

    for f in range(num_folds):
        X_test = X_train_folds[f]
        y_test = y_train_folds[f]
        X_train = np.concatenate(X_train_folds[:f] + X_train_folds[(f+
1):])
        y_train = np.concatenate(y_train_folds[:f] + y_train_folds[(f+
1):])

        knn.train(X_train, y_train)
        dists = knn.compute_distances(X_test, norm)
        y_test_pred = knn.predict_labels(dists, k = 10)

        num_correct = np.sum(y_test_pred == y_test)
        error = 1 - (num_correct / X_test.shape[0])
        errorsum += error

    avg_error = errorsum/num_folds
```

```

    Error_avg.append(avg_error)
    print ('Average error: %s (%.2f seconds)' % (avg_error, time.time(
)-t))

print('Total time: %.2f seconds' % (time.time()-time_start))

norms_name = ['L1_norm', 'L2_norm', 'Linf_norm']
plt.plot(norms_name, Error_avg, label = 'k-fold cross validation error
')
plt.xlabel('Norm used')
plt.ylabel('Cross-validation Error')
plt.grid()
plt.show()

```

**pass**

```

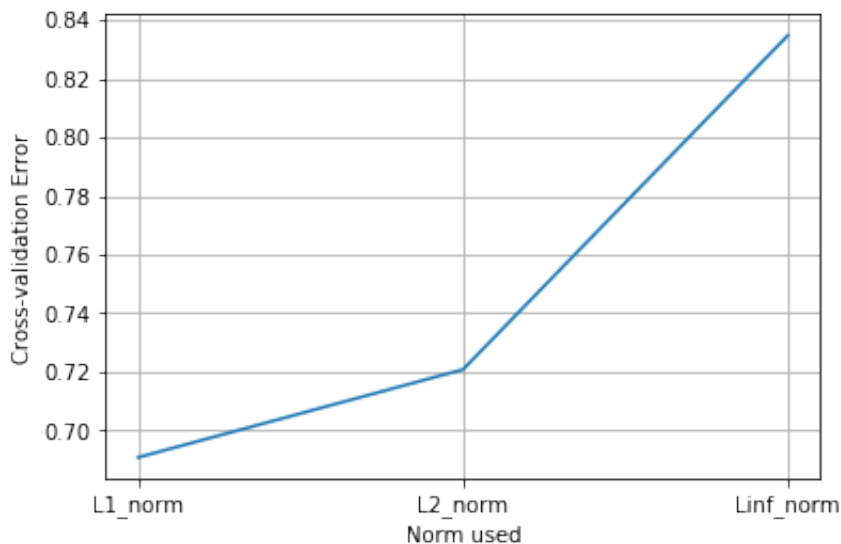
# ===== #
# END YOUR CODE HERE
# ===== #
print('Computation time: %.2f'%(time.time()-time_start))

```

```

Average error: 0.6908 (365.02 seconds)
Average error: 0.7208 (286.98 seconds)
Average error: 0.8348000000000001 (369.30 seconds)
Total time: 1021.30 seconds

```



Computation time: 1021.63

## 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) Cross-validation error is 0.6908 for the given norm and  $k = 10$ .

## 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 [327]: # ===== #
# YOUR CODE HERE:
# Evaluate the testing error of the k-nearest neighbors classifier
# for your optimal hyperparameters found by 5-fold cross-validation.
# ===== #
error = 1
optimal_k = 10
knn.train(X_train, y_train)
print (X_train.shape)
print (X_test.shape)

dists = knn.compute_distances(X = X_test, norm = L1_norm)
print(dists.shape)
y_test_pred = knn.predict_labels(dists, k = optimal_k)

error -= (np.sum(y_test_pred == y_test))/X_test.shape[0]

pass

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

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

(5000, 3072)
(500, 3072)
(500, 5000)
Error rate achieved: 0.722
```

## Question:

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

**Answer:**

Improved by around 0.01 in my case.

```

import numpy as np
import pdb

"""
This code was based off of code from cs231n at Stanford University, and
modified for ece239as 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)--data
        - y is a numpy array of size (num_examples, )--label
        """
        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].
                # ===== #
                dists[i,j] = norm(X[i]-self.X_train[j])
                pass
                # ===== #
                # 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.
    # ===== #

    Xtr = np.sum(self.X_train**2, axis = 1)          # (5000, 3072) -->
    # (5000, ^2) --> (5000, sum) --> (5000,)
    Xte = np.sum(X**2, axis = 1).reshape(X.shape[0], 1) # (500, 3072) -->
    # (500, ^2) --> (500, sum) --> (500,) --> (500, 1)
    dists = np.sqrt(Xtr + Xte - 2*X.dot(self.X_train.T)) # (5000,)+(500,1) ---
    # >(500, 5000)
    #pass

    # ===== #
    # 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.
    # ===== #
    neighbors_index = np.argsort(dists[i,:], axis = 0)
    #print (neighbors_index)
    closest_y = self.y_train[neighbors_index[:k]]
    #print(closest_y) #(500,k) with labels
    freq = np.bincount(closest_y)
    y_pred[i] = np.argmax(freq)

    # ===== #
    # END YOUR CODE HERE
    # ===== #
#print(len(y_pred))
return y_pred
```



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

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

Please print out the workbook entirely when completed.

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

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

## Importing libraries and data setup

```
In [2]: 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
```

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

```
%reload_ext autoreload
```

```
In [3]: # Set the path to the CIFAR-10 data
cifar10_dir = 'cifar-10-batches-py'
X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

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

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

```
In [4]: # 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', 'horse', 'ship', 'truck']
num_classes = len(classes)
samples_per_class = 7
for y, cls in enumerate(classes):
    idxs = np.flatnonzero(y_train == y)
    idxs = np.random.choice(idxs, samples_per_class, replace=False)
    for i, idx in enumerate(idxs):
        plt_idx = i * num_classes + y + 1
        plt.subplot(samples_per_class, num_classes, plt_idx)
        plt.imshow(X_train[idx].astype('uint8'))
        plt.axis('off')
        if i == 0:
            plt.title(cls)
plt.show()
```



```

In [5]: # Split the data into train, val, and test sets. In addition we will
# create a small development set as a subset of the training data;
# we can use this for development so our code runs faster.
num_training = 49000
num_validation = 1000
num_test = 1000
num_dev = 500

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

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

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

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

print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
print('Dev data shape: ', X_dev.shape)
print('Dev labels shape: ', y_dev.shape)

Train data shape: (49000, 32, 32, 3)
Train labels shape: (49000,)
Validation data shape: (1000, 32, 32, 3)
Validation labels shape: (1000,)
Test data shape: (1000, 32, 32, 3)
Test labels shape: (1000,)
Dev data shape: (500, 32, 32, 3)
Dev labels shape: (500,)

```

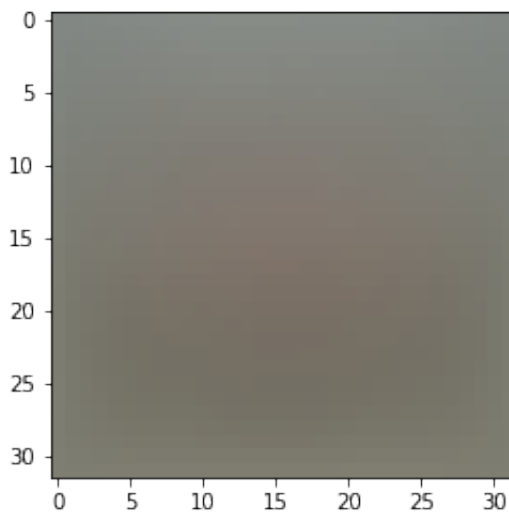
```
In [6]: # Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

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

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

```
In [7]: # 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]
```



```
In [8]: # second: subtract the mean image from train and test data
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image
```

```
In [9]: # 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)
```

## Question:

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

## Answer:

(1) For SVM, the mean-subtraction on the data serves to center the data since we will need to do gradient descent for training. Centering the data features allows the gradient descent to converge faster. For KNN, such step is not necessary since we don't need to do gradient descent. Also mean-subtraction in KNN will not change the distance that is used for calculating the nearest neighbor.

## Training an SVM

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

```
In [10]: from nndl.svm import SVM
```

```
In [11]: # Declare an instance of the SVM class.
# Weights are initialized to a random value.
# Note, to keep people's initial solutions consistent, we are going to use a random seed.

np.random.seed(1)

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

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

## SVM loss

```
In [12]: ## 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 15569.977915410193.
```

## SVM gradient

```
In [13]: ## 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: -9.406599 analytic: -9.406598, relative error: 2.569113e-  
08  
numerical: 1.794914 analytic: 1.794914, relative error: 8.310091e-08  
numerical: -3.751345 analytic: -3.751345, relative error: 1.407580e-  
08  
numerical: 8.445964 analytic: 8.445964, relative error: 4.708791e-08  
numerical: 5.482712 analytic: 5.482712, relative error: 2.485352e-08  
numerical: 14.075437 analytic: 14.075437, relative error: 2.081204e-  
08  
numerical: 4.225050 analytic: 4.225050, relative error: 2.273292e-08  
numerical: -13.966167 analytic: -13.966168, relative error: 1.765454  
e-08  
numerical: -2.306948 analytic: -2.306947, relative error: 1.098219e-  
07  
numerical: -11.544083 analytic: -11.544083, relative error: 2.504677  
e-08
```

## A vectorized version of SVM

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

```
In [14]: import time
```

```
In [15]: ## 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 {}'.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 {}'.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: 16059.146770808382 / 2073.798543122435 computed in 0.058722734451293945s
Vectorized loss / grad: 16059.146770808424 / 2073.798543122435 computed in 0.003857851028442383s
difference in loss / grad: -4.18367562815547e-11 / 2.8721770240751574e-12
```

## Stochastic gradient descent

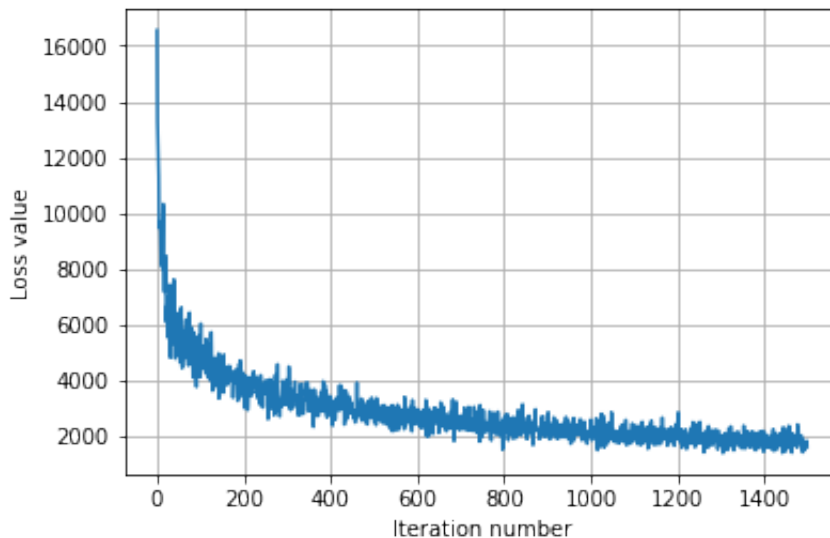
We now implement stochastic gradient descent. This uses the same principles of gradient descent we discussed in class, however, it calculates the gradient by only using examples from a subset of the training set (so each gradient calculation is faster).

```
In [16]: # Implement svm.train() by filling in the code to extract a batch of d
          # 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 {}'.format(toc - tic))

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

iteration 0 / 1500: loss 16557.380001909158
iteration 100 / 1500: loss 4701.089451272714
iteration 200 / 1500: loss 4017.333137942789
iteration 300 / 1500: loss 3681.9226471953625
iteration 400 / 1500: loss 2732.616437398899
iteration 500 / 1500: loss 2786.637842464506
iteration 600 / 1500: loss 2837.0357842782664
iteration 700 / 1500: loss 2206.2348687399326
iteration 800 / 1500: loss 2269.03882411698
iteration 900 / 1500: loss 2543.23781538592
iteration 1000 / 1500: loss 2566.692135726826
iteration 1100 / 1500: loss 2182.0689059051633
iteration 1200 / 1500: loss 1861.1182244250442
iteration 1300 / 1500: loss 1982.9013858528256
iteration 1400 / 1500: loss 1927.5204158582114
That took 5.759091138839722s
```



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



```
In [17]: ## Implement svm.predict() and use it to compute the training and test  
ing 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_p  
red)), ))  
  
training accuracy: 0.28530612244897957  
validation accuracy: 0.3
```

## Optimize the SVM

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

```

In [19]: # ===== #
# 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 VALIDAT
ION error.
#
#   Select the SVM that achieved the best validation error and report
#   its error rate on the test set.
#   Note: You do not need to modify SVM class for this section
# ===== #
l_r = [1e-10, 1e-9, 1e-8, 1e-7, 5e-7, 5e-6, 5e-5, 1e-5, 5e-4, 1e-3, 1e-
2]
best_valid_acc = 0
best_l_r = 0

for i in l_r:
    svm.train(X_train, y_train, learning_rate=i, num_iters=1500, batch
_size = 200, verbose=False)
    y_train_pred = svm.predict(X_train)
    y_val_pred = svm.predict(X_val)
    training_accuracy = (np.mean(np.equal(y_train, y_train_pred)))
    validation_accuracy = (np.mean(np.equal(y_val, y_val_pred)))
    if validation_accuracy > best_valid_acc:
        best_valid_acc = validation_accuracy
        best_l_r = i

print ('best validation accuracy', best_valid_acc, 'best validation er
ror =', float(1-best_valid_acc), 'learning rate,', best_l_r)

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

```

```

best validation accuracy 0.304 best validation error = 0.696 learnin
g rate, 0.001

```

```

import numpy as np
import pdb

"""
This code was based off of code from cs231n at Stanford University, and
modified for ece239as at UCLA.
"""
class SVM(object):

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

    def init_weights(self, dims):
        """
        Initializes the weight matrix of the SVM. 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)

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

        # compute the loss and the gradient
        num_classes = self.W.shape[0]
        num_train = X.shape[0]
        loss = 0.0
        #print(self.W.shape)

        # ===== #
        # YOUR CODE HERE:
        # Calculate the normalized SVM loss, and store it as '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.)
        # ===== #
        for i in np.arange(num_train):
            for j in np.arange(num_classes):
                if j == y[i]:
                    continue
                aj = self.W[j].dot(X[i])
                ay = self.W[y[i]].dot(X[i])
                loss += max(0, 1 + aj - ay)
        loss = loss / num_train

```

```

# ===== #
# 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.
    """
    # compute the loss and the gradient
    num_classes = self.W.shape[0]
    num_train = X.shape[0]
    loss = 0.0
    delta = 1
    grad = np.zeros_like(self.W)

    # ===== #
    # YOUR CODE HERE:
    # Calculate the SVM loss and the gradient. Store the gradient in
    # the variable grad.ears
    # ===== #
    for i in np.arange(num_train):
        for j in np.arange(num_classes):
            if j == y[i]:
                continue
            aj = self.W[j].dot(X[i]) #incorrect class score
            ay = self.W[y[i]].dot(X[i]) #correct class score
            z = delta+aj-ay
            loss += max(0, z)
            if z > 0:
                grad[j,:] +=X[i]
                grad[y[i],:] -=X[i]

    # ===== #
    # END YOUR CODE HERE
    # ===== #
    #print (grad.shape) #(10,3073)
    loss /= num_train
    grad /= num_train

    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.
    """
    num_classes = self.W.shape[0]
    num_train = X.shape[0]
    delta = 1
    loss = 0.0
    grad = np.zeros(self.W.shape) # initialize the gradient as zero

    # ===== #
    # YOUR CODE HERE:
    # Calculate the SVM loss WITHOUT any for loops.
    # ===== #
    scores = np.matmul(self.W,(X.T))
    #print('score shape', scores.shape) #(10,500)
    correct_scores = scores[y,np.arange(num_train)]
    #print('correct score shape', correct_scores.shape) #(500,)
    z = delta+scores-correct_scores
    margins = np.maximum(0, z) #broadcasting
    #print(margins.shape) (10,500)
    margins[y, np.arange(num_train)] = 0 #correct class score rows should be 0
    loss = np.sum(margins)
    loss = loss / num_train
    # ===== #
    # END YOUR CODE HERE
    # ===== #
    # YOUR CODE HERE:
    # Calculate the SVM grad WITHOUT any for loops.
    # ===== #
    grad_scores = np.zeros(margins.shape)
    grad_scores[margins > 0] = 1 #indicator function gives 1 if z>0, gives 0 if
    <=0
    #print(grad_scores.shape) #(10,500)
    grad_scores[y, np.arange(num_train)] = -np.sum(margins > 0, axis = 0)
    #print (X.shape) (500,3073)
    grad = np.matmul(grad_scores,(X))
    grad = grad / num_train
    # ===== #
    # 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: (batch_size, dim)
        # - 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)
        #print('indices shape', indices.shape) #(200,)
        X_batch = X[indices,:]
        #print('X_batch shape', X_batch.shape) #(200,3073)
        y_batch = y[indices]
        #print('y_batch shape', y_batch.shape) #(200,)

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

```

```

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

# ===== #
# YOUR CODE HERE:
#   Update the parameters, self.W, with a gradient step
# ===== #
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 with the parameter self.W.
    # ===== #
    scores = X.dot(self.W.T)
    #print(scores.shape)
    y_pred = np.argmax(scores, axis = 1)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return y_pred

```

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

Please follow the notebook linearly to implement a softmax classifier.

Please print out the workbook entirely when completed.

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

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

In [2]:

```
import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

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

In [3]:

```
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'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # subsample the data
    mask = list(range(num_training, num_training + num_validation))
    X_val = X_train[mask]
    y_val = y_train[mask]
    mask = list(range(num_training))
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = list(range(num_test))
    X_test = X_test[mask]
    y_test = y_test[mask]
    mask = np.random.choice(num_training, num_dev, replace=False)
    X_dev = X_train[mask]
```



```

y_dev = y_train[mask]

# Preprocessing: reshape the image data into rows
X_train = np.reshape(X_train, (X_train.shape[0], -1))
X_val = np.reshape(X_val, (X_val.shape[0], -1))
X_test = np.reshape(X_test, (X_test.shape[0], -1))
X_dev = np.reshape(X_dev, (X_dev.shape[0], -1))

# Normalize the data: subtract the mean image
mean_image = np.mean(X_train, axis = 0)
X_train -= mean_image
X_val -= mean_image
X_test -= mean_image
X_dev -= mean_image

# add bias dimension and transform into columns
X_train = np.hstack([X_train, np.ones((X_train.shape[0], 1))])
X_val = np.hstack([X_val, np.ones((X_val.shape[0], 1))])
X_test = np.hstack([X_test, np.ones((X_test.shape[0], 1))])
X_dev = np.hstack([X_dev, np.ones((X_dev.shape[0], 1))])

return X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test, X_dev, y_dev = get_CIFAR10_data(
)
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)
print('dev data shape: ', X_dev.shape)
print('dev labels shape: ', y_dev.shape)

```

```

Train data shape: (49000, 3073)
Train labels shape: (49000,)
Validation data shape: (1000, 3073)
Validation labels shape: (1000,)
Test data shape: (1000, 3073)
Test labels shape: (1000,)
dev data shape: (500, 3073)
dev labels shape: (500,)

```

## Training a softmax classifier.

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

In [4]:

```
from nndl import Softmax
```

In [5]:

```
# Declare an instance of the Softmax class.  
# Weights are initialized to a random value.  
# Note, to keep people's first solutions consistent, we are going to use a random seed.  
  
np.random.seed(1)  
  
num_classes = len(np.unique(y_train))  
num_features = X_train.shape[1]  
  
softmax = Softmax(dims=[num_classes, num_features])
```

## Softmax loss

In [6]:

```
## Implement the loss function of the softmax using a for loop over  
# the number of examples  
  
loss = softmax.loss(X_train, y_train)
```

In [7]:

```
print(loss)
```

2.327760702804897

## Question:

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

## Answer:

In a properly configured dataset, the expected initial loss should be about  $\log(\text{num\_classes})$ . Here number of classes is 10, and  $\log(10)$  is about 2.3, thus it makes sense.

## Softmax gradient

In [8]:

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

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

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

```
numerical: 1.654102 analytic: 1.654102, relative error: 7.045946e-09
numerical: 0.478814 analytic: 0.478814, relative error: 1.117901e-07
numerical: 2.294308 analytic: 2.294308, relative error: 9.050479e-09
numerical: 1.078868 analytic: 1.078868, relative error: 4.821854e-08
numerical: 0.644949 analytic: 0.644949, relative error: 1.443827e-07
numerical: 0.552875 analytic: 0.552875, relative error: 3.957662e-08
numerical: 0.825816 analytic: 0.825816, relative error: 6.564272e-08
numerical: -1.645459 analytic: -1.645459, relative error: 2.527257e-
08
numerical: 0.021223 analytic: 0.021223, relative error: 6.634177e-07
numerical: -1.591346 analytic: -1.591346, relative error: 4.446770e-
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 [9]:

```
import time
```

In [10]:

```
## 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.
```

```
Normal loss / grad_norm: 2.3115405774531843 / 297.680496535171 computed in 0.07746696472167969s
Vectorized loss / grad: 2.311540577453182 / 297.680496535171 computed in 0.004712343215942383s
difference in loss / grad: 2.220446049250313e-15 / 1.9612927145041015e-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?

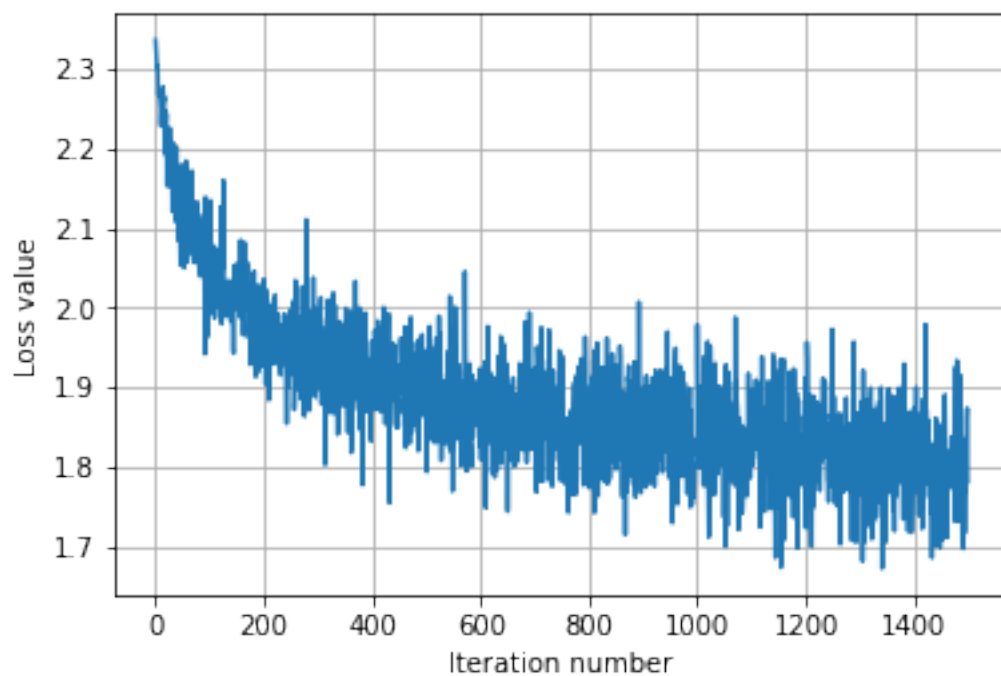
## Answer:

The gradient descent in both cases are the same, following the same equation. However, the losses and gradients are defined differently, where the softmax computes probabilities of the labels, and SVM outputs uncalibrated scores. In terms of learning rate, the softmax should have a smaller learning rate given that we take exponentials for the gradient. Consequently, the learning rate of softmax is much smaller than that of SVM to achieve the same variation in weights.

In [11]:

```
# 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.grid()  
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 1.8992158530357477
iteration 1000 / 1500: loss 1.9783503540252299
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 6.10860800743103s
```



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

In [12]:

```
## 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
```

# Optimize the softmax classifier

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

In [13]:

```
np.finfo(float).eps
```

Out[13]:

```
2.220446049250313e-16
```

In [14]:

```
# ===== #
# YOUR CODE HERE:
#   Train the Softmax classifier with different learning rates and
#   evaluate on the validation data.
#   Report:
#     - The best learning rate of the ones you tested.
#     - The best validation accuracy corresponding to the best validation error.
#
#   Select the SVM that achieved the best validation error and report
#   its error rate on the test set.
# ===== #
l_r = [1e-10, 1e-9, 1e-8, 1e-7, 5e-7, 5e-6, 5e-5]
best_valid_acc = 0
best_l_r = 0

for i in l_r:
    softmax.train(X_train, y_train, learning_rate=i, num_iters=1500, batch_size
= 200, verbose=False)
    y_train_pred = softmax.predict(X_train)
    y_val_pred = softmax.predict(X_val)
    training_accuracy = (np.mean(np.equal(y_train, y_train_pred)))
    validation_accuracy = (np.mean(np.equal(y_val, y_val_pred)))
    if validation_accuracy > best_valid_acc:
        best_valid_acc = validation_accuracy
        best_l_r = i

print ('best validation accuracy', best_valid_acc, 'Error rate', 1-best_valid_ac
c, 'learning rate,', best_l_r)

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

```
best validation accuracy 0.404 Error rate 0.596 learning rate, 5e-07
```

```

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.)
        # ===== #
        num_classes = self.W.shape[0]
        num_train = X.shape[0]

        for i in np.arange(num_train):
            a = self.W.dot(X[i].T) #(10,)
            a -= np.max(a) #normalization
            sum_as = np.sum(np.exp(a))
            normalize = np.exp(a) / sum_as
            #print(np.exp(a))
            loss += -np.log(normalize[y[i]]) #exp(a_i)/sum
            #for j in np.arange(num_classes):

        loss = loss / num_train

        # ===== #
        # 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.
    # ===== #
    num_classes = self.W.shape[0]
    num_train = X.shape[0]

    for i in np.arange(num_train):
        a = self.W.dot(X[i].T) #(10,)
        a -= np.max(a) #normalization
        sum_as = np.sum(np.exp(a))
        softmaxes = np.exp(a) / sum_as
        #print(np.exp(a))
        loss += -np.log(softmaxes[y[i]]) #exp(a_i)/sum

        for j in np.arange(num_classes):
            if j == y[i]:
                grad[j,:] += (-1 + softmaxes[y[i]])*(X[i])
                #grad[j,:] += (-1 + softmaxes[j])*(X[i])
            else:
                grad[j,:] += softmaxes[j]*(X[i])

    #print('grad shape', grad.shape)
    loss = loss / num_train
    grad = grad / num_train
    # ===== #
    # 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 ouptuts as loss_and_grad.
    """
    loss = 0.0
    grad = np.zeros(self.W.shape) # initialize the gradient as zero
    #print (X.shape) (500,3073)
    # ===== #
    # YOUR CODE HERE:
    # Calculate the softmax loss and gradient WITHOUT any for loops.
    # ===== #
    num_classes = self.W.shape[0]
    num_train = X.shape[0]

    a = np.matmul(self.W,(X.T)) #(C,N) (10,500)
    a -= np.max(a, axis = 0, keepdims = True) #shift a so that the highest
        number is 0, max for every sample
    #print(np.max(a, axis = 0, keepdims = True).shape) #(10,1)

    #sum_a = np.sum(np.exp(a), axis = 0, keepdims = True) #(10,1)
    softmaxes = np.exp(a) / np.sum(np.exp(a), axis = 0, keepdims = True)
        #(10,500) is this broadcasting too?

    loss = np.sum(-np.log(softmaxes[y, np.arange(num_train)]))

    grad_scores = np.zeros(softmaxes.shape) #np.zeros_like
    grad_scores[y, np.arange(num_train)] = 1 #indicator function gives 1 if y =
        j, otherwise 0 #(10,500)
    grad = (softmaxes-grad_scores).dot(X)

    #print('grad shape', grad.shape)

    loss = loss / num_train
    grad = grad / num_train

    # ===== #
    # 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: (batch_size, dim)
        # - 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)
        #print('indices shape', indices.shape) #(200,)
        X_batch = X[indices,:]
        #print('X_batch shape', X_batch.shape) #(200,3073)
        y_batch = y[indices]
        #print('y_batch shape', y_batch.shape) #(200,)
        # ===== #
        # END YOUR CODE HERE
        # ===== #

```

```

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

# ===== #
# YOUR CODE HERE:
#   Update the parameters, self.W, with a gradient step
# ===== #
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[0])
    # ===== #
    # YOUR CODE HERE:
    #   Predict the labels given the training data.
    # ===== #
    scores = X.dot(self.W.T)
    #print(scores.shape)
    y_pred = np.argmax(scores, axis = 1)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    return y_pred

```