

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

## Assignment #2

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

Please print out the workbook entirely when completed.

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

## Import the appropriate libraries

```
In [1]: import numpy as np # for doing most of our calculations
import matplotlib.pyplot as plt # for plotting
from utils.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
```

```
In [2]: # 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)
```

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

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



In [4]:

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

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

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

(5000, 3072) (500, 3072)

## K-nearest neighbors

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

In [5]:

```
# Import the KNN class

from nndl import KNN
```

In [6]:

```
# Declare an instance of the knn class.
knn = KNN()

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

## Questions

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

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

# Answers

(1) In the function `knn.train()`, we simply assign our training data and labels into the knn model.

(2) Pros: It is simple and fast, only  $O(1)$  time complexity.

Cons: 1. It is memory-intensive because we need to store all the input data. If there are huge amounts of data, it will cost lots of memory. 2. Making prediction will be inefficient, and it will cost much time on computation.

## KNN prediction

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

In [7]:

```
# Implement the function compute_distances() in the KNN class.
# Do not worry about the input 'norm' for now; use the default definition 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: 34.11346912384033
Frobenius norm of L2 distances: 7906696.077040902
```

### Really slow code

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

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

## KNN vectorization

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

In [8]:

```
# 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): {}'.for
```

```
Time to run code: 0.19804787635803223
Difference in L2 distances between your KNN implementations (should be 0): 0.0
```

### Speedup

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

## Implementing the prediction

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

In [9]:

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

error = 1

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

y_pred = knn.predict_labels(dists_L2_vectorized, 1)

num_incorrect = 0
for i in np.arange(num_test):
    if y_pred[i] != y_test[i]:
        num_incorrect += 1

error = num_incorrect/num_test

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

print(error)
```

0.726

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

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

## Optimizing KNN hyperparameters

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

### Create training and validation folds

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

In [10]:

```
# Create the dataset folds for cross-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]
# ===== #

X_train_folds = np.split(X_train, num_folds)
y_train_folds = np.split(y_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 [11]:

```

time_start = time.time()

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

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

errors = []
for k in ks:
    error = 0
    for i in np.arange(num_folds):
        # Train
        knn = KNN()
        X_train_fold = np.concatenate(X_train_folds[:i] + X_train_folds[(i+1):])
        y_train_fold = np.concatenate(y_train_folds[:i] + y_train_folds[(i+1):])
        X_test_fold = X_train_folds[i]
        y_test_fold = y_train_folds[i]
        knn.train(X=X_train_fold, y=y_train_fold)

        # Distance
        dists_L2_vectorized = knn.compute_L2_distances_vectorized(X_test_fold)

        # Predict
        y_pred = knn.predict_labels(dists_L2_vectorized, k)
        num_incorrect = 0
        num_test_fold = y_test_fold.shape[0]
        for j in np.arange(num_test_fold):
            if y_pred[j] != y_test_fold[j]:
                num_incorrect += 1
        error += num_incorrect/num_test_fold

    errors.append(error/num_folds)

print("Errors: {}".format(errors))
plt.scatter(ks, errors)
plt.title("k vs. Cross-Validation Error")
plt.xlabel("ks")

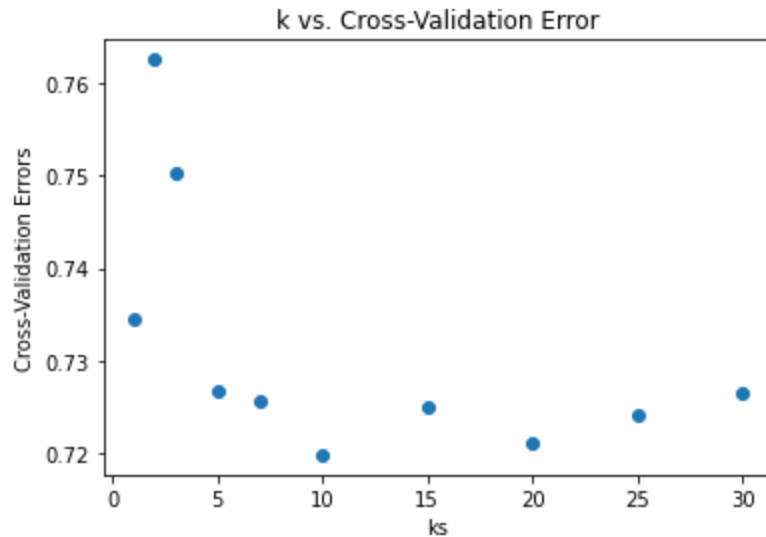
```

```
plt.ylabel("Cross-Validation Errors")
plt.show()
```

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

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

Errors: [0.7344, 0.7626000000000002, 0.7504000000000001, 0.7267999999999999, 0.7256, 0.7198, 0.725, 0.721, 0.7242, 0.7266]



Computation time: 25.91

## 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)  $k = 10$
- (2) error = 0.7198

## Optimizing the norm

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

In [12]:

```
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.
# ===== #

k = 10
errors = []
for norm in norms:
    error = 0
    for i in np.arange(num_folds):
        # Train
        knn = KNN()
        X_train_fold = np.concatenate(X_train_folds[:i] + X_train_folds[(i+1):])
        y_train_fold = np.concatenate(y_train_folds[:i] + y_train_folds[(i+1):])
        X_test_fold = X_train_folds[i]
        y_test_fold = y_train_folds[i]
        knn.train(X=X_train_fold, y=y_train_fold)

        # Distance
        dists = knn.compute_distances(X=X_test_fold, norm=norm)

        # Predict
        y_pred = knn.predict_labels(dists, k)
        num_incorrect = 0
        num_test_fold = y_test_fold.shape[0]
        for j in np.arange(num_test_fold):
            if y_pred[j] != y_test_fold[j]:
                num_incorrect += 1
        error += num_incorrect/num_test_fold

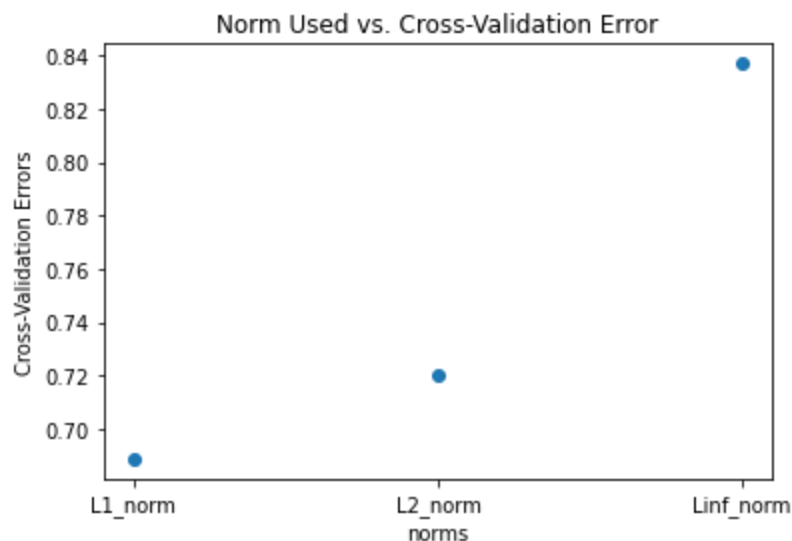
    errors.append(error/num_folds)

print("Errors: {}".format(errors))
norms_names = ['L1_norm', 'L2_norm', 'Linf_norm']
plt.scatter(norms_names, errors)
plt.title("Norm Used vs. Cross-Validation Error")
plt.xlabel("norms")
plt.ylabel("Cross-Validation Errors")
plt.show()

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

```

Errors: [0.6886000000000001, 0.7198, 0.8370000000000001]



## 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
- (2) error = 0.6886

## Evaluating the model on the testing dataset.

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

In [13]:

```
error = 1

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

k = 10
num_test = 500
# Train
knn = KNN()
knn.train(X=X_train, y=y_train)

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

# Predict
y_pred = knn.predict_labels(dists, k)
num_incorrect = 0
for i in np.arange(num_test):
    if y_pred[i] != y_test[i]:
        num_incorrect += 1

error = num_incorrect/num_test

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

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

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:

k=1 and using the L2-norm: 0.726

k=10 and using the L1-norm: 0.722

Improvement:  $0.726 - 0.722 = 0.004$

```

1 import numpy as np
2 import pdb
3
4
5 class KNN(object):
6
7     def __init__(self):
8         pass
9
10    def train(self, X, y):
11        """
12        Inputs:
13        - X is a numpy array of size (num_examples, D)
14        - y is a numpy array of size (num_examples, )
15        """
16        self.X_train = X
17        self.y_train = y
18
19    def compute_distances(self, X, norm=None):
20        """
21        Compute the distance between each test point in X and each training point
22        in self.X_train.
23
24        Inputs:
25        - X: A numpy array of shape (num_test, D) containing test data.
26        - norm: the function with which the norm is taken.
27
28        Returns:
29        - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
30            is the Euclidean distance between the ith test point and the jth
31            training point.
32        """
33        if norm is None:
34            norm = lambda x: np.sqrt(np.sum(x**2))
35            #norm = 2
36
37        num_test = X.shape[0]
38        num_train = self.X_train.shape[0]
39        dists = np.zeros((num_test, num_train))
40
41        for i in np.arange(num_test):
42            for j in np.arange(num_train):
43                # ===== #
44                # YOUR CODE HERE:
45                #   Compute the distance between the ith test point and the jth
46                #   training point using norm(), and store the result in dists[i, j].
47
48                # ===== #
49
50                dist = norm(X[i] - self.X_train[j])
51                dists[i][j] = dist
52
53                # ===== #
54                # END YOUR CODE HERE
55                # ===== #
56
57        return dists

```

```

58 def compute_L2_distances_vectorized(self, X):
59     """
60     Compute the distance between each test point in X and each training point
61     in self.X_train WITHOUT using any for loops.
62
63     Inputs:
64     - X: A numpy array of shape (num_test, D) containing test data.
65
66     Returns:
67     - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
68         is the Euclidean distance between the ith test point and the jth
69         training point.
70     """
71     num_test = X.shape[0]
72     num_train = self.X_train.shape[0]
73     dists = np.zeros((num_test, num_train))
74
75     # ===== #
76     # YOUR CODE HERE:
77     # Compute the L2 distance between the ith test point and the jth
78     # training point and store the result in dists[i, j]. You may
79     # NOT use a for loop (or list comprehension). You may only use
80     # numpy operations.
81     #
82     # HINT: use broadcasting. If you have a shape (N,1) array and
83     # a shape (M,) array, adding them together produces a shape (N, M)
84     # array.
85     # ===== #
86
87     X2 = np.sum(X**2, axis=1).reshape((num_test, 1)) # shape is (num_test, 1)
88     Y2 = np.sum(self.X_train**2, axis=1).reshape((1, num_train)) # shape is
89     (1, num_train)
90     XY = X.dot(self.X_train.T) # shape is (num_test, num_train)
91     dists = np.sqrt(X2 + Y2 - 2*XY) # shape is (num_test, num_train)
92
93     # ===== #
94     # END YOUR CODE HERE
95     # ===== #
96
97     return dists
98
99 def predict_labels(self, dists, k=1):
100     """
101     Given a matrix of distances between test points and training points,
102     predict a label for each test point.
103
104     Inputs:
105     - dists: A numpy array of shape (num_test, num_train) where dists[i, j]
106         gives the distance between the ith test point and the jth training
107         point.
108
109     Returns:
110     - y: A numpy array of shape (num_test,) containing predicted labels for
111         the test data, where y[i] is the predicted label for the test point X[i].
112     """
113     num_test = dists.shape[0]
114     y_pred = np.zeros(num_test)

```

```

114 for i in np.arange(num_test):
115     # A list of length k storing the labels of the k nearest neighbors to
116     # the ith test point.
117     closest_y = []
118     # ===== #
119     # YOUR CODE HERE:
120     # Use the distances to calculate and then store the labels of
121     # the k-nearest neighbors to the ith test point. The function
122     # numpy.argsort may be useful.
123     #
124     # After doing this, find the most common label of the k-nearest
125     # neighbors. Store the predicted label of the ith training example
126     # as y_pred[i]. Break ties by choosing the smaller label.
127     # ===== #
128
129     sortedIdxs = np.argsort(dists[i])
130     closest_y = self.y_train[sortedIdxs[:k]]
131     y_pred[i] = np.argmax(np.bincount(closest_y))
132
133     # ===== #
134     # END YOUR CODE HERE
135     # ===== #
136
137 return y_pred
138

```

2.

$$\text{Let } a_i(x) = w_i^T x + b_i, \quad \theta = \{w_i, b_i\},$$

$$i = 1, \dots, c$$

$$\begin{aligned} \mathcal{L}(\theta) &= \prod_{i=1}^m p(y^{(i)} | x^{(i)}, \theta) \\ \downarrow \text{likelihood} \\ &= \prod_{i=1}^m \text{softmax}_{y^{(i)}}(x^{(i)}) \\ &= \prod_{i=1}^m \frac{e^{a_{y^{(i)}}(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} \end{aligned}$$

$$\Rightarrow \log \mathcal{L}(\theta) = \sum_{i=1}^m \log \left[ \frac{e^{a_{y^{(i)}}(x^{(i)})}}{\sum_{j=1}^c e^{a_j(x^{(i)})}} \right]$$

$$= \sum_{i=1}^m \left[ a_{y^{(i)}}(x^{(i)}) - \log \left( \sum_{j=1}^c e^{a_j(x^{(i)})} \right) \right]$$

#

$$\Rightarrow \underset{\theta}{\operatorname{argmax}} \log \mathcal{L}(\theta) = \underset{\theta}{\operatorname{argmin}} \underbrace{-\log \mathcal{L}(\theta)}$$

↓ negative log-likelihood  
( $\mathcal{L}$ )  
we want in this  
question

Let  $J(\theta) = -\log L(\theta)$

$$\nabla_{w_i} J(\theta) = \sum_{k=1}^m \frac{e^{a_i(x^{(k)})}}{\sum_{j=1}^c e^{a_j(x^{(k)})}} x^{(k)} - \sum_{\{k: y^{(k)}=i, \forall k \in \{1, \dots, m\}\}} x^{(k)}$$

$$= \sum_{k=1}^m \text{softmax}_i(x^{(k)}) x^{(k)} - \sum_{\{k: y^{(k)}=i, \forall k \in \{1, \dots, m\}\}} x^{(k)}$$

$$\nabla_{b_i} J(\theta) = \sum_{k=1}^m \frac{e^{a_i(x^{(k)})}}{\sum_{j=1}^c e^{a_j(x^{(k)})}} - \sum_{\{k: y^{(k)}=i, \forall k \in \{1, \dots, m\}\}} 1$$

$$= \sum_{k=1}^m \text{softmax}_i(x^{(k)}) - \sum_{\{k: y^{(k)}=i, \forall k \in \{1, \dots, m\}\}} 1$$

#

# 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.

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

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

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

```
In [2]: 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]
        mask = list(range(num_test))
        X_test = X_test[mask]
        y_test = y_test[mask]
        mask = np.random.choice(num_training, num_dev, replace=False)
        X_dev = X_train[mask]
        y_dev = y_train[mask]

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

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

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

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

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

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

## Training a softmax classifier.

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

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

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

np.random.seed(1)

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

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

### Softmax loss

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

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

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

```
2.3277607028048966
```

## Question:

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



# Answer:

It is because we have 10 classes and it is randomly distribution. We can expect the probability of each class is 0.1. Therefore,  $-\log(0.1) \approx 2.3$  is reasonable to be the loss.

## Softmax gradient

In [7]:

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

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

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

```
numerical: -0.597248 analytic: -0.597248, relative error: 1.352238e-09
numerical: 1.231072 analytic: 1.231072, relative error: 5.063644e-09
numerical: -1.063401 analytic: -1.063401, relative error: 3.143855e-09
numerical: 1.924655 analytic: 1.924655, relative error: 1.512317e-08
numerical: 1.005965 analytic: 1.005965, relative error: 5.340822e-08
numerical: 2.265621 analytic: 2.265621, relative error: 1.803286e-08
numerical: -1.235339 analytic: -1.235339, relative error: 3.233409e-08
numerical: -1.706993 analytic: -1.706993, relative error: 8.622539e-09
numerical: -1.196225 analytic: -1.196225, relative error: 8.455186e-09
numerical: -2.421361 analytic: -2.421361, relative error: 2.693823e-08
```

## A vectorized version of Softmax

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

In [8]:

```
import time
```

In [9]:

```
## 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), 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), 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: 2.325139432933735 / 364.46996613435 computed in 0.0562129020690918s
```

Vectorized loss / grad: 2.3251394329337334 / 364.46996613435 computed in 0.0017008781433105469s  
difference in loss / grad: 1.7763568394002505e-15 / 3.048411952454789e-13

## Stochastic gradient descent

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

### Question:

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

### Answer:

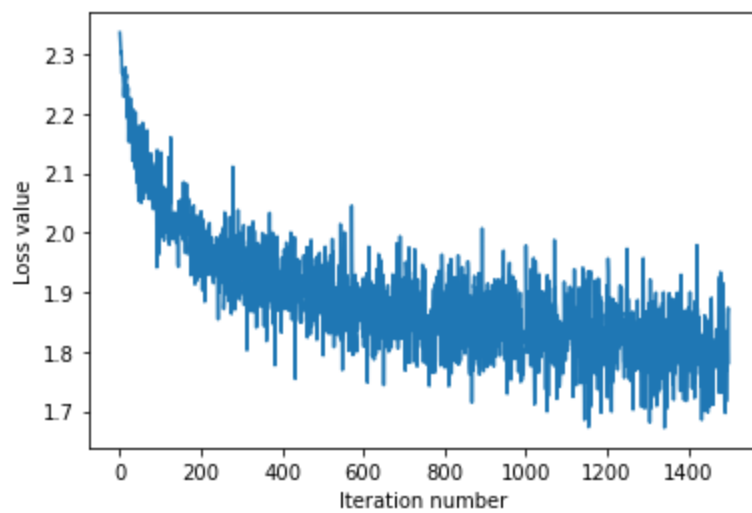
They are identical. The processes have no difference.

```
In [10]: # 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 2.0357745120662813
iteration 300 / 1500: loss 1.9813348165609888
iteration 400 / 1500: loss 1.9583142443981614
iteration 500 / 1500: loss 1.862265307354135
iteration 600 / 1500: loss 1.8532611454359382
iteration 700 / 1500: loss 1.835306222372583
iteration 800 / 1500: loss 1.8293892468827635
iteration 900 / 1500: loss 1.8992158530357484
iteration 1000 / 1500: loss 1.97835035402523
iteration 1100 / 1500: loss 1.8470797913532633
iteration 1200 / 1500: loss 1.8411450268664082
iteration 1300 / 1500: loss 1.79104024957921
iteration 1400 / 1500: loss 1.8705803029382257
That took 3.599605083465576s
```



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

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

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

training accuracy: 0.3811428571428571
validation accuracy: 0.398
```

## Optimize the softmax classifier

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

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

```
Out[12]: 2.220446049250313e-16
```

```
In [13]: # ===== #
# 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.
# ===== #

learning_rates = [10**i for i in range(-9, 0)]
accuracy = {}

best_learning_rate = 0
best_validation = 0

for learning_rate in learning_rates:
    softmax = Softmax(dims=[num_classes, num_features])
    loss = softmax.train(X_train, y_train, learning_rate=learning_rate, num_iters=1500, ve
```

```

y_train_pred = softmax.predict(X_train)
train_accuracy = np.mean(np.equal(y_train, y_train_pred))

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

accuracy[learning_rate] = (train_accuracy, val_accuracy)

if best_validation < val_accuracy:
    best_learning_rate = learning_rate
    best_validation = val_accuracy

for learning_rate in accuracy:
    print("Learning Rate: {}, Train Accuracy: {}, Validation: {}".format(learning_rate, ac

print("\nThe Best Learning Rate: {}".format(best_learning_rate))
print("The Best Validation Accuracy: {}".format(best_validation))
print("The Best Validation Error: {}\n".format(1 - best_validation))

# Best Test
softmax.train(X_train, y_train, learning_rate=best_learning_rate, num_iters=1500, verbose=
y_test_pred = softmax.predict(X_test)

test_accuracy = np.mean(np.equal(y_test, y_test_pred))
print("Test Accuracy: {}\nError rate on the test set: {}".format(test_accuracy, 1-test_acc

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

```

/Users/jacky/My\_Data/Data/UCLA/2022\_Winter/ECE C247\_Deep Learning/Homework/HW2/hw2-code/nn  
dl/softmax.py:142: RuntimeWarning: divide by zero encountered in log

```

probs_log = -np.log(probs_row)
Learning Rate: 1e-09, Train Accuracy: 0.17079591836734695, Validation: 0.16
Learning Rate: 1e-08, Train Accuracy: 0.2886938775510204, Validation: 0.304
Learning Rate: 1e-07, Train Accuracy: 0.38210204081632654, Validation: 0.395
Learning Rate: 1e-06, Train Accuracy: 0.42248979591836733, Validation: 0.407
Learning Rate: 1e-05, Train Accuracy: 0.34916326530612246, Validation: 0.33
Learning Rate: 0.0001, Train Accuracy: 0.2804081632653061, Validation: 0.264
Learning Rate: 0.001, Train Accuracy: 0.2738979591836735, Validation: 0.255
Learning Rate: 0.01, Train Accuracy: 0.2778775510204082, Validation: 0.27
Learning Rate: 0.1, Train Accuracy: 0.2919387755102041, Validation: 0.284

```

The Best Learning Rate: 1e-06  
The Best Validation Accuracy: 0.407  
The Best Validation Error: 0.593

Test Accuracy: 0.402  
Error rate on the test set: 0.598

```

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

```

```

57     return loss
58
59
60 def loss_and_grad(self, X, y):
61     """
62     Same as self.loss(X, y), except that it also returns the gradient.
63
64     Output: grad -- a matrix of the same dimensions as W containing
65             the gradient of the loss with respect to W.
66     """
67
68     # Initialize the loss and gradient to zero.
69     loss = 0.0
70     grad = np.zeros_like(self.W)
71
72     # ===== #
73     # YOUR CODE HERE:
74     # Calculate the softmax loss and the gradient. Store the gradient
75     # as the variable grad.
76     # ===== #
77
78     a = self.W.dot(X.T).T
79
80     i = 0
81     for row in a:
82         row -= np.max(row) #avoid overflow
83         a_row = np.sum(np.exp(row))
84         loss += (np.log(a_row) - row[y[i]])
85
86         for j in np.arange(self.W.shape[0]):
87             grad[j] += (np.exp(row[j])/a_row) * X[i]
88             grad[y[i]] -= X[i]
89             i += 1
90
91     loss /= a.shape[0]
92     grad /= a.shape[0]
93
94     # ===== #
95     # END YOUR CODE HERE
96     # ===== #
97
98     return loss, grad
99
100 def grad_check_sparse(self, X, y, your_grad, num_checks=10, h=1e-5):
101     """
102     sample a few random elements and only return numerical
103     in these dimensions.
104     """
105
106     for i in np.arange(num_checks):
107         ix = tuple([np.random.randint(m) for m in self.W.shape])
108
109         oldval = self.W[ix]
110         self.W[ix] = oldval + h # increment by h
111         fxph = self.loss(X, y)
112         self.W[ix] = oldval - h # decrement by h
113         fxmh = self.loss(X,y) # evaluate f(x - h)
114         self.W[ix] = oldval # reset
115
116         grad_numerical = (fxph - fxmh) / (2 * h)

```

```

117     grad_analytic = your_grad[ix]
118     rel_error = abs(grad_numerical - grad_analytic) / (abs(grad_numerical)
+ abs(grad_analytic))
119     print('numerical: %f analytic: %f, relative error: %e' %
(grad_numerical, grad_analytic, rel_error))
120
121 def fast_loss_and_grad(self, X, y):
122     """
123     A vectorized implementation of loss_and_grad. It shares the same
124     inputs and outputs as loss_and_grad.
125     """
126     loss = 0.0
127     grad = np.zeros(self.W.shape) # initialize the gradient as zero
128
129     # ===== #
130     # YOUR CODE HERE:
131     # Calculate the softmax loss and gradient WITHOUT any for loops.
132     # ===== #
133
134     a = self.W.dot(X.T).T
135     num_train = a.shape[0]
136
137     a -= np.max(a, axis=1, keepdims=True)
138     a_exp = np.exp(a)
139
140     probs = a_exp / np.sum(a_exp, axis=1, keepdims=True)
141     probs_row = probs[range(num_train), y]
142     probs_log = -np.log(probs_row)
143
144     loss = np.sum(probs_log) / num_train
145
146     probs[range(num_train), y] -= 1
147     grad = (probs.T.dot(X)) / num_train
148
149     # ===== #
150     # END YOUR CODE HERE
151     # ===== #
152
153     return loss, grad
154
155 def train(self, X, y, learning_rate=1e-3, num_iters=100,
156         batch_size=200, verbose=False):
157     """
158     Train this linear classifier using stochastic gradient descent.
159
160     Inputs:
161     - X: A numpy array of shape (N, D) containing training data; there are N
162         training samples each of dimension D.
163     - y: A numpy array of shape (N,) containing training labels; y[i] = c
164         means that X[i] has label 0 ≤ c < C for C classes.
165     - learning_rate: (float) learning rate for optimization.
166     - num_iters: (integer) number of steps to take when optimizing
167     - batch_size: (integer) number of training examples to use at each step.
168     - verbose: (boolean) If true, print progress during optimization.
169
170     Outputs:
171     A list containing the value of the loss function at each training
172     iteration.
173     """
174     num_train, dim = X.shape

```

```

174     num_classes = np.max(y) + 1 # assume y takes values 0...K-1 where K is
number of classes
175
176     self.init_weights(dims=[np.max(y) + 1, X.shape[1]]) # initializes the
weights of self.W
177
178     # Run stochastic gradient descent to optimize W
179     loss_history = []
180
181     for it in np.arange(num_iters):
182         X_batch = None
183         y_batch = None
184
185         # ===== #
186         # YOUR CODE HERE:
187         #     Sample batch_size elements from the training data for use in
188         #     gradient descent. After sampling,
189         #     - X_batch should have shape: (dim, batch_size)
190         #     - y_batch should have shape: (batch_size,)
191         #     The indices should be randomly generated to reduce correlations
192         #     in the dataset. Use np.random.choice. It's okay to sample with
193         #     replacement.
194         # ===== #
195
196         indices = np.random.choice(X.shape[0], batch_size)
197         X_batch = X[indices]
198         y_batch = y[indices]
199
200         # ===== #
201         # END YOUR CODE HERE
202         # ===== #
203
204         # evaluate loss and gradient
205         loss, grad = self.fast_loss_and_grad(X_batch, y_batch)
206         loss_history.append(loss)
207
208         # ===== #
209         # YOUR CODE HERE:
210         #     Update the parameters, self.W, with a gradient step
211         # ===== #
212
213         self.W -= learning_rate * grad
214
215         # ===== #
216         # END YOUR CODE HERE
217         # ===== #
218
219         if verbose and it % 100 == 0:
220             print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
221
222     return loss_history
223
224 def predict(self, X):
225     """
226     Inputs:
227     - X: N x D array of training data. Each row is a D-dimensional point.
228
229     Returns:
230     - y_pred: Predicted labels for the data in X. y_pred is a 1-dimensional
231       array of length N, and each element is an integer giving the predicted

```



```
232     class.  
233     """  
234     y_pred = np.zeros(X.shape[1])  
235     # ===== #  
236     # YOUR CODE HERE:  
237     #   Predict the labels given the training data.  
238     # ===== #  
239  
240     a = self.W.dot(X.T).T  
241     y_pred = np.argmax(a, axis=1)  
242  
243     # ===== #  
244     # END YOUR CODE HERE  
245     # ===== #  
246  
247     return y_pred  
248  
249
```