

## 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 [164]: 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
```

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

```
In [165]: # 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 [166]: # 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 [167]: # 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 [168]: # Import the KNN class

from nn1 import KNN
```

```
In [169]: # Declare an instance of the knn class.
knn = KNN()

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

## Questions

- (1) Describe what is going on in the function `knn.train()`.
- (2) What are the pros and cons of this training step?

## Answers

- (1) The function `knn.train()` simply remembers the data points (images) and their labels.
- (2) The pros are that there is no training time and it is simple and fast; the cons are that it is memory intensive as we must store all the input data. Additionally, predicting the test data is expensive since we must get the distance to all points.

## 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 [170]: # 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: 19.605123281478882  
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 [171]: # Implement the function compute_L2_distances_vectorized() in the KNN class.
# In this function, you ought to achieve the same L2 distance but WITHOUT any for loops.
# Note, this is SPECIFIC for the L2 norm.

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

Time to run code: 0.14534711837768555  
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 [172]: # 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, k=1)
num_incorrect = 0
for i in range(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 [173]: # Create the dataset folds for cross-validation.
num_folds = 5

X_train_folds = []
y_train_folds = []

np.random.seed(123)
idx = np.random.permutation(num_training)
print(idx, idx.shape)
# ===== #
# 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_shuffle = X_train[idx]
y_train_shuffle = y_train[idx]
X_train_folds = np.array_split(X_train_shuffle, num_folds)
y_train_folds = np.array_split(y_train_shuffle, num_folds)

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

print(y_train_folds[0].shape)
```

```
[2648 2456 4557 ... 1346 3454 3582] (5000,)
(1000,)
```

### 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 [174]: 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.
# ===== #
error_list = []
for k in ks:
    error = 0
    for i in np.arange(num_folds):
        X_fold_train = np.concatenate(X_train_folds[:i] + X_train_folds[i + 1:], axis=0)
        y_fold_train = np.concatenate(y_train_folds[:i] + y_train_folds[i + 1:], axis=0)
        X_fold_val = X_train_folds[i]
        y_fold_val = y_train_folds[i]

        knn.train(X_fold_train, y_fold_train)
        cur_pred = knn.predict_labels(knn.compute_L2_distances_vectorized(X_fold_val), k)

        num_incorrect = 0
        for j in range(len(cur_pred)):
            if cur_pred[j] != y_fold_val[j]:
                num_incorrect += 1
        error += num_incorrect / X_fold_val.shape[0]
    error_list.append(error / num_folds)

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

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

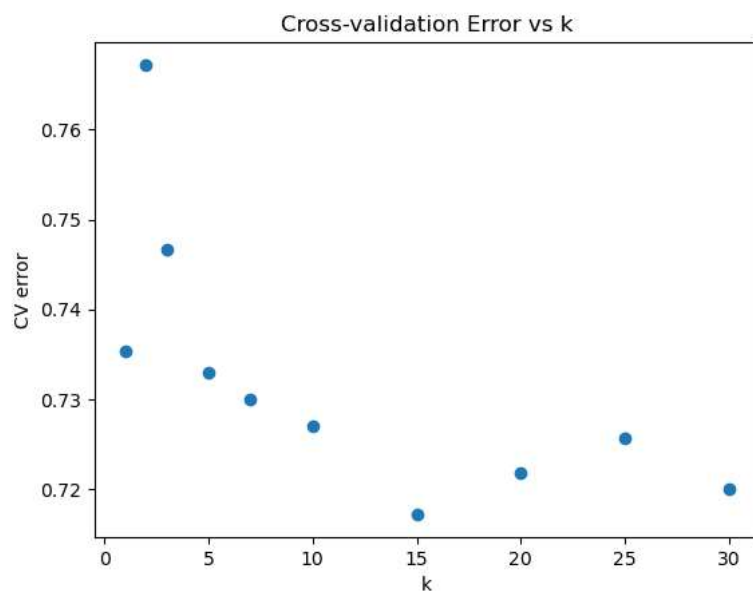
```

Computation time: 20.47

```
In [175]: plt.scatter(ks, error_list)
plt.title('Cross-validation Error vs k')
plt.ylabel('CV error')
plt.xlabel('k')
plt.show()

for i in range(len(error_list)):
    print(str(ks[i]) + ': ' + str(error_list[i]))

print()
best_error = min(error_list)
k_idx = error_list.index(best_error)
best_k = ks[k_idx]
print("Best k, error: " + str(best_k) + ', ' + str(best_error))
```



```
1: 0.7354
2: 0.7672
3: 0.7465999999999999
5: 0.733
7: 0.73
10: 0.727
15: 0.7172
20: 0.7218
25: 0.7259999999999999
30: 0.72
```

```
Best k, error: 15, 0.7172
```

## 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)  $k=15$
- (2) 0.7172

## 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 [179]: 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.
# ===== #
norm_error_list = []
for norm in norms:
    norm_error = 0
    for i in np.arange(num_folds):
        X_fold_train = np.concatenate(X_train_folds[:i] + X_train_folds[i + 1:], axis=0)
        y_fold_train = np.concatenate(y_train_folds[:i] + y_train_folds[i + 1:], axis=0)
        X_fold_val = X_train_folds[i]
        y_fold_val = y_train_folds[i]

        knn.train(X_fold_train, y_fold_train)
        cur_pred = knn.predict_labels(knn.compute_distances(X_fold_val, norm), best_k)

        num_incorrect = 0
        for j in range(len(cur_pred)):
            if cur_pred[j] != y_fold_val[j]:
                num_incorrect += 1
        norm_error += num_incorrect / X_fold_val.shape[0]
    norm_error_list.append(norm_error / num_folds)

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

```

Computation time: 354.56

```

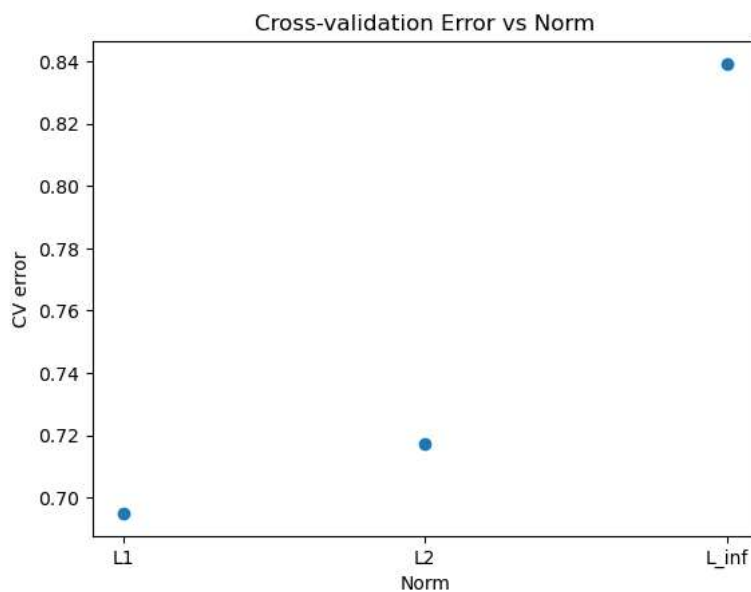
In [180]: ns = ["L1", "L2", "L_inf"]

plt.scatter(["L1", "L2", "L_inf"], norm_error_list)
plt.title('Cross-validation Error vs Norm')
plt.ylabel('CV error')
plt.xlabel('Norm')
plt.show()

for i in range(len(norm_error_list)):
    print(ns[i] + ' error: ' + str(norm_error_list[i]))

print()
best_error = min(norm_error_list)
norm_idx = norm_error_list.index(best_error)
print("Best norm, error: " + ns[norm_idx] + ', ' + str(best_error))

```



L1 error: 0.695  
 L2 error: 0.7172  
 L\_inf error: 0.8392

Best norm, error: L1, 0.695

## 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) 0.695

## 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 [181]: error = 1
# ===== #
# YOUR CODE HERE:
# Evaluate the testing error of the k-nearest neighbors classifier
# for your optimal hyperparameters found by 5-fold cross-validation.
# ===== #

knn = KNN()
knn.train(X=X_train, y=y_train)
dists_L1 = knn.compute_distances(X=X_test, norm=L1_norm)
y_pred = knn.predict_labels(dists_L1, k=best_k)

num_incorrect = 0
for i in range(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: 0.718

## Question:

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

## Answer:

Error with L2-norm and  $k=1$ : 0.726

Error with L1-norm and  $k=15$ : 0.718

Improvement: 0.008

```

In [ ]: import numpy as np
import pdb

class KNN(object):

    def __init__(self):
        pass

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

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

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

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

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

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

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

                # ===== #
                # 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.
        # ===== #
        # print(X.shape) # (500, 3072)
        # print(self.X_train.shape) # (5000, 3072)

```

```

# X**2.sum(axis=1).reshape(-1,1): (N,1)
# X_train**2.sum(axis=1): (M, )
# X @ X_train.T: (500, 3072) x (3072, 5000) = (N, M)
dists = np.sqrt((X**2).sum(axis=1).reshape(-1,1) - 2 * X @ self.X_train.T + (self.X_train**2).sum(axis=1))

# ===== #
# 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.
        # ===== #
        indices = np.argsort(dists[i])
        closest_y = self.y_train[indices][:k]
        closest_y = sorted(list(closest_y)) # sort by label number, then get most common label
        y_pred[i] = max(closest_y, key = closest_y.count)

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

    return y_pred

```