

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

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

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class ([cs231n.stanford.edu](http://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 [91]: 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 data
from collections import Counter
import operator

# 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-ipy
%load_ext autoreload
%autoreload 2
```

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

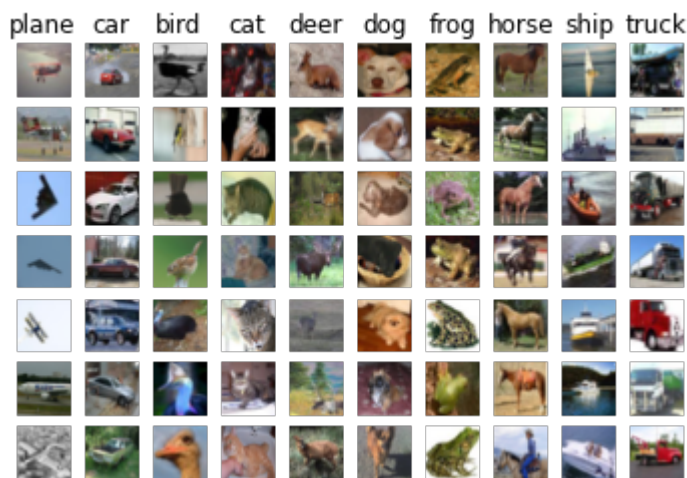
```
%reload_ext autoreload
```

```
In [80]: # Set the path to the CIFAR-10 data
cifar10_dir = '/home/alon/school/c247a/datasets/cifar-10-batches-py' # You need
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 [81]: # 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', 'sh
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 [82]: # 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 [8]: `# Import the KNN class`

```
from nndl import KNN
```

In [9]: `# 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 program is caching the entire training dataset, in order to be able to compute distances using the chosen metric for each test example.

(2) Pros: this step is very fast and simple,  $O(1)$  training complexity. Cons: if the dataset is large, this takes a lot of memory. Further, the amount of memory needed scales with the number of training examples and the input vector size.

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

```
print('Shape of dists_L2: {}'.format(dists_L2.shape))
```

```
Time to run code: 35.81522274017334
```

```
Frobenius norm of L2 distances: 7906696.077040902
```

```
Shape of dists_L2: (500, 5000)
```

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 [153... # 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  
# 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)  
0  
Time to run code: 0.23537254333496094  
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 [185...

```

k = 1

correct = 0

num_test = dists_L2_vectorized.shape[0]
y_pred = np.zeros(num_test)

for i in np.arange(num_test):
    sorted_row_k_ixs = np.argsort(dists_L2_vectorized[i])[:k]
    #sorted_row_k = np.sort(dists_L2_vectorized[i])[:k]
    #print(sorted_row_k_ixs)
    nn_class_ixs = [y_train[ix] for ix in sorted_row_k_ixs]
    #print(nn_class_ixs)
    nn_class_top = max(set(nn_class_ixs), key=nn_class_ixs.count)
    #print(nn_class_top)
    #print("predicted class: {}".format(classes[nn_class_top]))
    #print("correct class: {}".format(classes[y_test[i]]))

    y_pred[i] = nn_class_top

    if (nn_class_top == y_test[i]):
        correct += 1
    #print('\n')

print("Total correct: {}".format(correct))
print("Percent correct: {}".format(correct/num_test))

```

Total correct: 137  
Percent correct: 0.274

In [131...

```

# 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
num_errors = 0

# ===== #
# 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.
# ===== #

predicted_labels = knn.predict_labels(dists_L2_vectorized, k=1)
#print(predicted_labels)
for i, lab in enumerate(predicted_labels):
    if (lab != y_test[i]):
        num_errors += 1

error = num_errors/y_test.shape[0]

#pass
# ===== #
# 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 [136... # 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]
# ===== #

fold_size = int(X_train.shape[0]/num_folds)
#print(fold_size)

for i in range(num_folds):
    X_fold = X_train[fold_size*i:fold_size*(i+1)]
    X_train_folds.append(X_fold)
    y_fold = y_train[fold_size*i:fold_size*(i+1)]
    y_train_folds.append(y_fold)

# ===== #
# 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 [150...

```
#sad = [int(item) for item in np.arange(len(X_train_folds))!=2]  
sad = [x for i,x in enumerate(X_train_folds) if i!=3]  
bad = np.concatenate(sad,axis=0)  
print(bad.shape)
```

(4000, 3072)

```

In [176... 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.
# ===== #

errs = []

#loop through the ks:
for k in ks:

    errs_k = []

    #loop through the folds:
    for i in range(0, num_folds):
        #obtain validation fold
        X_val_fold = X_train_folds[i]
        y_val_fold = y_train_folds[i]

        #obtain training fold set
        X_train_set = np.concatenate([x for j,x in enumerate(X_train_folds) if j != i])
        y_train_set = np.concatenate([y for j,y in enumerate(y_train_folds) if j != i])

        # Declare an instance of the knn class.
        knn_inst = 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_inst.train(X=X_train_set, y=y_train_set)

        dists = knn_inst.compute_L2_distances_vectorized(X=X_val_fold)
        predicted_labels = knn_inst.predict_labels(dists, k=k)

        num_errors = 0

        for index, label in enumerate(predicted_labels):
            if label != y_val_fold[index]:
                num_errors += 1

        #print("Number of errors for this step: {}".format(num_errors))
        error = num_errors/y_val_fold.shape[0]

        #print("Cross-validation error for fold {} and k={}: {}".format(i, k, error))
        errs_k.append(error)

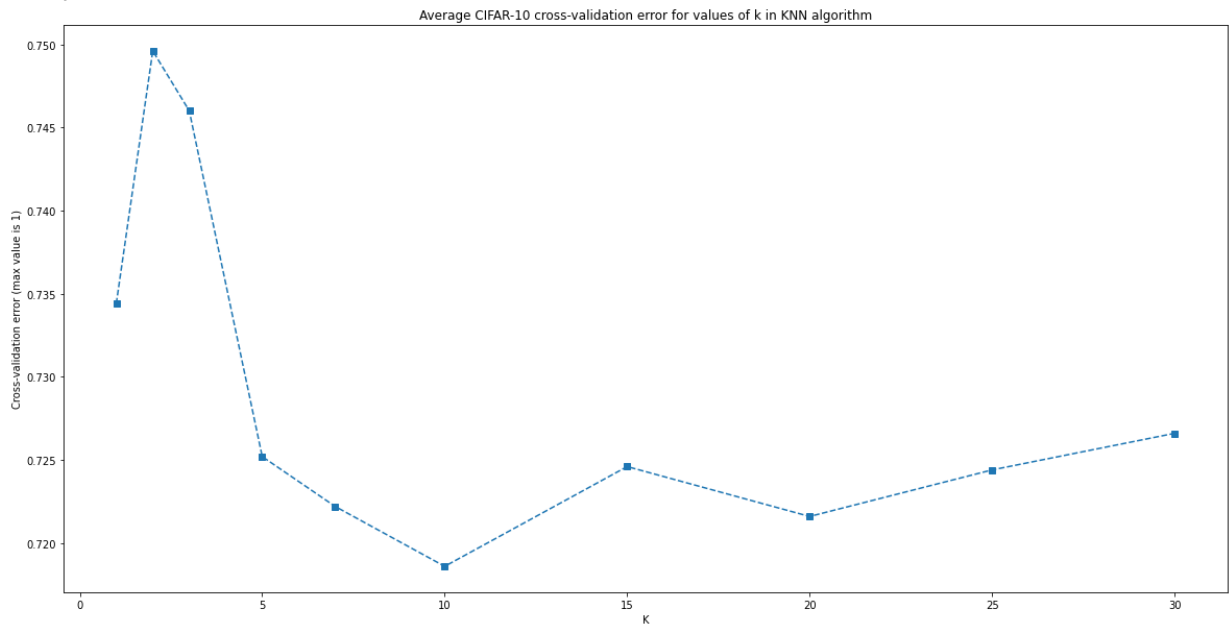
    errs.append(errs_k)

#average errors for each k to plot
average_errs_k = [sum(lst)/len(lst) for lst in errs]
plt.figure(figsize=(20,10))

```



Computation time: 35.57



## 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)  $k=10$  has the best performance.
- (2) 0.7186 or 71.86%

## 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 [181... 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.
# ===== #

errs = []

for norm in norms:
    #print(norm())

    errs_norm = []

    #loop through the folds:
    for i in range(0, num_folds):
        #obtain validation fold
        X_val_fold = X_train_folds[i]
        y_val_fold = y_train_folds[i]

        #obtain training fold set
        X_train_set = np.concatenate([x for j,x in enumerate(X_train_folds) if j!=i])
        y_train_set = np.concatenate([y for j,y in enumerate(y_train_folds) if j!=i])

        # Declare an instance of the knn class.
        knn_inst = 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_inst.train(X=X_train_set, y=y_train_set)

        #dists = knn_inst.compute_L2_distances_vectorized(X=X_val_fold)
        dists = knn_inst.compute_distances(X=X_val_fold, norm=norm)
        #print(dists.shape)
        predicted_labels = knn_inst.predict_labels(dists, k=10)

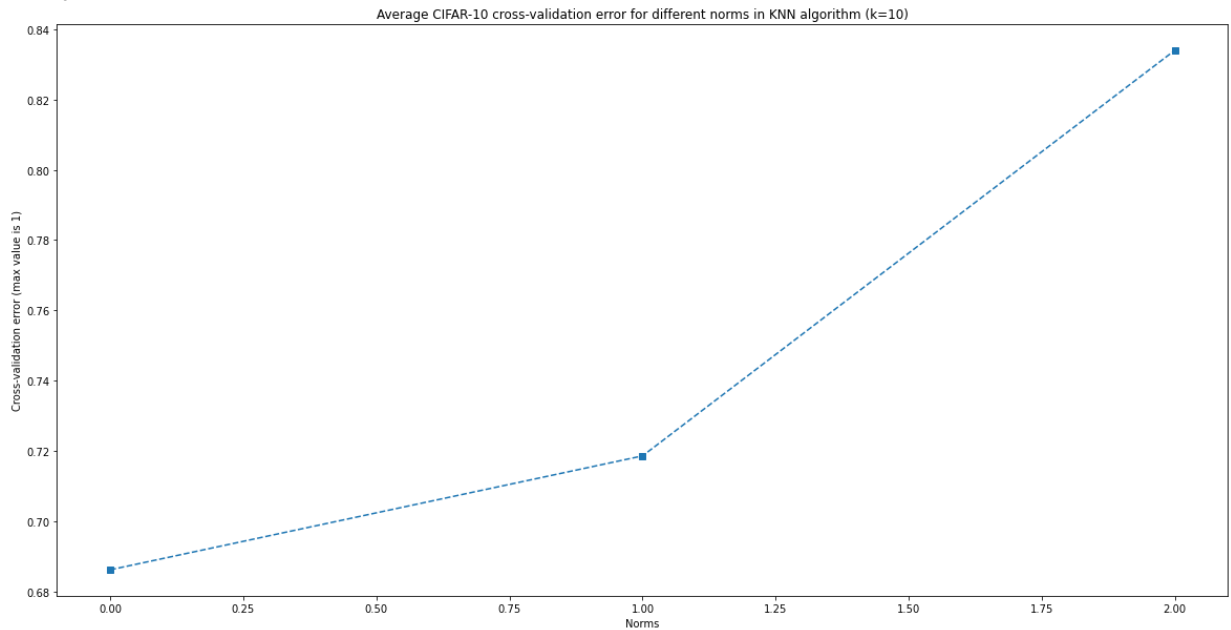
        num_errors = 0

        for index, label in enumerate(predicted_labels):
            if label != y_val_fold[index]:
                num_errors += 1

        #print("Number of errors for this step: {}".format(num_errors))
        error = num_errors/y_val_fold.shape[0]

```

Computation time: 738.26



## Questions:

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

In [182...

```
print(average_errs_norm)
```

```
[0.6862000000000001, 0.7186, 0.834]
```

## Answers:

- (1) The L1 norm.
- (2) 0.6862 or 68.62%

## 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 [184...

```

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.
# ===== #

#using k = 10 and L1 norm:

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

#train
knn_inst.train(X=X_train, y=y_train)

#compute distances and predict labels
dists = knn_inst.compute_distances(X=X_test, norm=L1_norm)
predicted_labels = knn_inst.predict_labels(dists, k=10)

num_errors = 0

for index, label in enumerate(predicted_labels):
    if label != y_test[index]:
        num_errors += 1

#print("Number of errors for this step: {}".format(num_errors))
error = num_errors/y_test.shape[0]

print("Error rate for k=10 and L1 Norm KNN algorithm: {}".format(error))
print("Improvement with cross validation approach: {}".format(0.726 - error))

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

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

```

```

Error rate for k=10 and L1 Norm KNN algorithm: 0.714
Improvement with cross validation approach: 0.012000000000000001
Error rate achieved: 0.714

```

## Question:

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

## Answer:

I got an improvement of 0.012 or 1.2%.

```

In [ ]: import numpy as np
import pdb
from collections import Counter
import operator

"""
This code was based off of code from cs231n at Stanford University, and modified
"""

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 j
                # training point using norm(), and store the result in dists[i, j].
                # ===== #

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

            #pass

```