

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]: 1 import numpy as np # for doing most of our calculations
        2 import matplotlib.pyplot as plt # for plotting
        3 from utils.data_utils import load_CIFAR10 # function to load the CIFAR-10 data
        4
        5 # Load matplotlib images inline
        6 %matplotlib inline
        7
        8 # These are important for reloading any code you write in external .py files
        9 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
       10 %load_ext autoreload
       11 %autoreload 2
```

```
In [2]: 1 # Set the path to the CIFAR-10 data
        2 cifar10_dir = 'dataset\cifar-10-batches-py' # You need to update this line
        3 X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
        4
        5 # As a sanity check, we print out the size of the training and test data.
        6 print('Training data shape: ', X_train.shape)
        7 print('Training labels shape: ', y_train.shape)
        8 print('Test data shape: ', X_test.shape)
        9 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]:

```

1  # Visualize some examples from the dataset.
2  # We show a few examples of training images from each class.
3  classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 's
4  num_classes = len(classes)
5  samples_per_class = 7
6  for y, cls in enumerate(classes):
7      idxs = np.flatnonzero(y_train == y)
8      idxs = np.random.choice(idxs, samples_per_class, replace=False)
9      for i, idx in enumerate(idxs):
10         plt_idx = i * num_classes + y + 1
11         plt.subplot(samples_per_class, num_classes, plt_idx)
12         plt.imshow(X_train[idx].astype('uint8'))
13         plt.axis('off')
14         if i == 0:
15             plt.title(cls)
16 plt.show()

```



```
In [4]: 1 # Subsample the data for more efficient code execution in this exercise
2 num_training = 5000
3 mask = list(range(num_training))
4 X_train = X_train[mask]
5 y_train = y_train[mask]
6
7 num_test = 500
8 mask = list(range(num_test))
9 X_test = X_test[mask]
10 y_test = y_test[mask]
11
12 # Reshape the image data into rows
13 X_train = np.reshape(X_train, (X_train.shape[0], -1))
14 X_test = np.reshape(X_test, (X_test.shape[0], -1))
15 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]: 1 # Import the KNN class
2
3 from nnd1 import KNN
```

```
In [6]: 1 # Declare an instance of the knn class.
2 knn = KNN()
3
4 # Train the classifier.
5 # We have implemented the training of the KNN classifier.
6 # Look at the train function in the KNN class to see what this does.
7 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 `knn.train()` function, we are simply storing all the training data samples and the corresponding labels in memory.
- (2) An obvious pro is that the training step is extremely fast and simple to implement. However, a con is memory usage. As the amount of samples increase, the model needs more storage.

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]: 1 # Implement the function compute_distances() in the KNN class.
2 # Do not worry about the input 'norm' for now; use the default definition of
3 # in the code, which is the 2-norm.
4 # You should only have to fill out the clearly marked sections.
5
6 import time
7 time_start =time.time()
8
9 dists_L2 = knn.compute_distances(X=X_test)
10
11 print('Time to run code: {}'.format(time.time()-time_start))
12 print('Frobenius norm of L2 distances: {}'.format(np.linalg.norm(dists_L2, 'fro')))
```

Time to run code: 47.018675088882446

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]: 1 # Implement the function compute_L2_distances_vectorized() in the KNN class.
2 # In this function, you ought to achieve the same L2 distance but WITHOUT an
3 # Note, this is SPECIFIC for the L2 norm.
4
5 time_start =time.time()
6 dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_test)
7 print('Time to run code: {}'.format(time.time()-time_start))
8 print('Difference in L2 distances between your KNN implementations (should be 0): 0.0')
```

Time to run code: 0.5874133110046387

Difference in L2 distances between your KNN implementations (should be 0): 0.0

```
In [12]: 1 print(dists_L2_vectorized.shape)
```

(500, 5000)

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]: 1 # Implement the function predict_labels in the KNN class.
2 # Calculate the training error (num_incorrect / total_samples)
3 # from running knn.predict_labels with k=1
4
5 error = 1
6
7 # ===== #
8 # YOUR CODE HERE:
9 # Calculate the error rate by calling predict_labels on the test
10 # data with k = 1. Store the error rate in the variable error.
11 # ===== #
12 distances = knn.compute_L2_distances_vectorized(X=X_test)
13 predictions = knn.predict_labels(distances,1)
14 error = sum(predictions!=y_test)/len(y_test)
15 # ===== #
16 # END YOUR CODE HERE
17 # ===== #
18
19 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]: 1 # Create the dataset folds for cross-validation.
2 num_folds = 5
3
4 X_train_folds = []
5 y_train_folds = []
6 np.random.seed(24)
7 # ===== #
8 # YOUR CODE HERE:
9 # Split the training data into num_folds (i.e., 5) folds.
10 # X_train_folds is a list, where X_train_folds[i] contains the
11 # data points in fold i.
12 # y_train_folds is also a list, where y_train_folds[i] contains
13 # the corresponding labels for the data in X_train_folds[i]
14 # ===== #
15 indices = np.arange(0,X_train.shape[0])
16 randomIndices=np.random.permutation(indices)
17 foldSize = int(X_train.shape[0]/num_folds)
18 start = 0
19 end = start + foldSize
20 for i in range(0,num_folds-1):
21     indices = randomIndices[start:end]
22     foldtraindata=X_train[indices]
23     foldtrainlabel = y_train[indices]
24     X_train_folds.append(foldtraindata)
25     y_train_folds.append(foldtrainlabel)
26     start = end
27     end = start+foldSize
28 indices = randomIndices[start:]
29 foldtraindata=X_train[indices]
30 foldtrainlabel = y_train[indices]
31 X_train_folds.append(foldtraindata)
32 y_train_folds.append(foldtrainlabel)
33
34 # ===== #
35 # END YOUR CODE HERE
36 # ===== #
37
38

```

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]: 1 time_start =time.time()
2
3 ks = [1, 2, 3, 5, 7, 10, 15, 20, 25, 30]
4
5 # ===== #
6 # YOUR CODE HERE:
7 # Calculate the cross-validation error for each k in ks, testing
8 # the trained model on each of the 5 folds. Average these errors
9 # together and make a plot of k vs. cross-validation error. Since
10 # we are assuming L2 distance here, please use the vectorized code!
11 # Otherwise, you might be waiting a long time.
12 # ===== #
13 knn = KNN()
14 kError = []
15 for k in ks:
16     foldError = []
17     for i in range(0,len(X_train_folds)):
18         trainCopy = X_train_folds.copy()
19         trainLabelsCopy = y_train_folds.copy()
20         foldTest = trainCopy[i]
21         foldTestlabel = trainLabelsCopy[i]
22         trainCopy.pop(i)
23         trainLabelsCopy.pop(i)
24         foldTr=np.array(trainCopy)
25         foldTrLab=np.array(trainLabelsCopy)
26         foldTr = foldTr.reshape([foldTr.shape[0]*foldTr.shape[1],foldTr.shap
27         foldTrLab = foldTrLab.reshape([foldTrLab.shape[0]*foldTrLab.shape[1]
28         knn.train(foldTr,foldTrLab)
29         distances = knn.compute_L2_distances_vectorized(X=foldTest)
30         predictions = knn.predict_labels(distances,k)
31         error = sum(predictions!=foldTestlabel)/len(foldTestlabel)
32         foldError.append(error)
33     print("Done cross validation for k %d. Error is %.2f"%(k,np.mean(foldErr
34     kError.append(np.mean(foldError))
35
36
37 plt.plot(ks, kError, 'bo--')
38 plt.axis([0, 31, 0.71, 0.77])
39 plt.xlabel('k value')
40 plt.ylabel('mean cross validation error')
41 plt.show()
42 # ===== #
43 # END YOUR CODE HERE
44 # ===== #
45
46 print('Computation time: %.2f'%(time.time()-time_start))

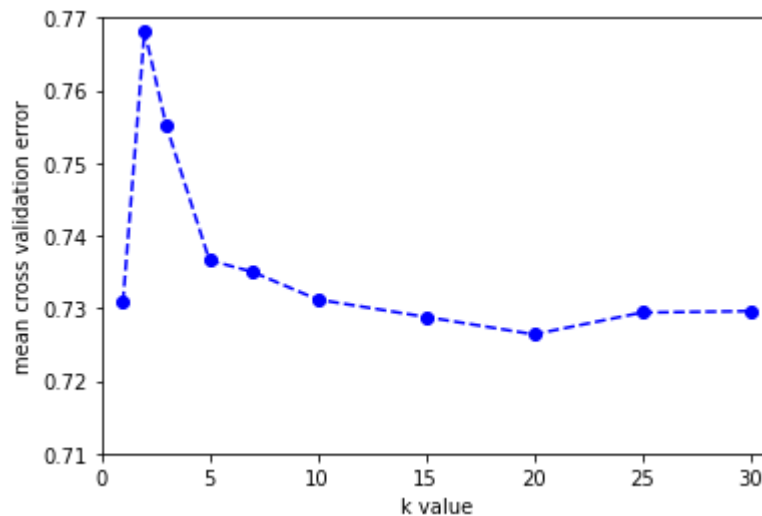
```

```

Done cross validation for k 1. Error is 0.73
Done cross validation for k 2. Error is 0.77
Done cross validation for k 3. Error is 0.76
Done cross validation for k 5. Error is 0.74
Done cross validation for k 7. Error is 0.73
Done cross validation for k 10. Error is 0.73
Done cross validation for k 15. Error is 0.73
Done cross validation for k 20. Error is 0.73

```

Done cross validation for k 25. Error is 0.73
Done cross validation for k 30. Error is 0.73



Computation time: 35.01

In [12]:

1	kError
---	--------

Out[12]: [0.7307999999999999,
0.7681999999999999,
0.7552000000000001,
0.7365999999999999,
0.735,
0.7312,
0.7288,
0.7264,
0.7293999999999999,
0.7295999999999999]

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) The best value of k is found to be 20.

(2) The cross validation error is 0.7264.

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 [13]: 1 time_start =time.time()
2
3 L1_norm = lambda x: np.linalg.norm(x, ord=1)
4 L2_norm = lambda x: np.linalg.norm(x, ord=2)
5 Linf_norm = lambda x: np.linalg.norm(x, ord= np.inf)
6 norms = [L1_norm, L2_norm, Linf_norm]
7 normNames=["l1","l2","inf"]
8 # ===== #
9 # YOUR CODE HERE:
10 # Calculate the cross-validation error for each norm in norms, testing
11 # the trained model on each of the 5 folds. Average these errors
12 # together and make a plot of the norm used vs the cross-validation error
13 # Use the best cross-validation k from the previous part.
14 #
15 # Feel free to use the compute_distances function. We're testing just
16 # three norms, but be advised that this could still take some time.
17 # You're welcome to write a vectorized form of the L1- and Linf- norms
18 # to speed this up, but it is not necessary.
19 # ===== #
20 bestK=20
21 knn = KNN()
22 normError = []
23 for ind in range(0,len(norms)):
24     normType = norms[ind]
25     normName = normNames[ind]
26     foldError = []
27     for i in range(0,len(X_train_folds)):
28         trainCopy = X_train_folds.copy()
29         trainLabelsCopy = y_train_folds.copy()
30         foldTest = trainCopy[i]
31         foldTtestlabel = trainLabelsCopy[i]
32         trainCopy.pop(i)
33         trainLabelsCopy.pop(i)
34         foldTr=np.array(trainCopy)
35         foldTrLab=np.array(trainLabelsCopy)
36         foldTr = foldTr.reshape([foldTr.shape[0]*foldTr.shape[1],foldTr.shap
37         foldTrLab = foldTrLab.reshape([foldTrLab.shape[0]*foldTrLab.shape[1]
38         knn.train(foldTr,foldTrLab)
39         if normName=="l2":
40             distances = knn.compute_L2_distances_vectorized(foldTest)
41         else:
42             distances = knn.compute_distances(foldTest,normType)
43
44         #distances = knn.compute_distances(X=foldTest,norm=normType)
45         predictions = knn.predict_labels(distances,k)
46         error = sum(predictions!=foldTtestlabel)/len(foldTtestlabel)
47         print("Fold done")
48         foldError.append(error)
49     print("Done cross validation for norm %. Error is %.2f"%(normName,np.me
50     normError.append(np.mean(foldError))
51 plt.figure()
52 normNames=['L1_norm', 'L2_norm', 'Linf_norm']
53 plt.plot(normNames,normError)
54 plt.xlabel('Norm Type')
55 plt.ylabel('Mean Cross Val Error')
56 # ===== #

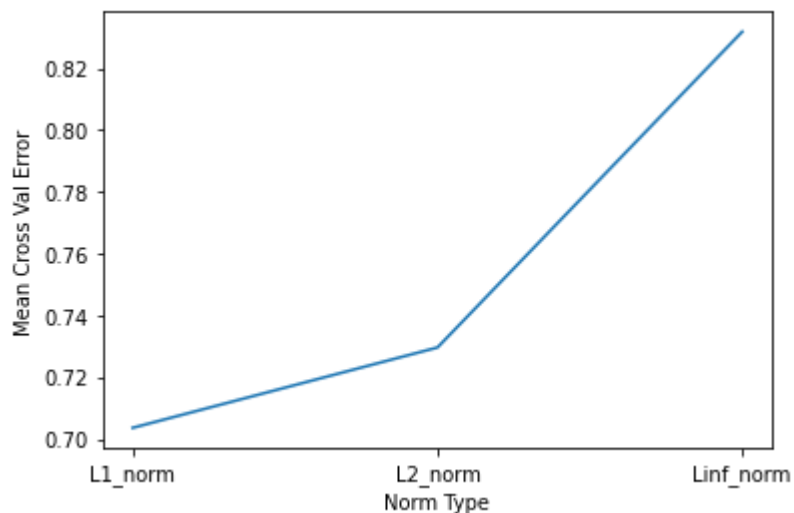
```

```

57 # END YOUR CODE HERE
58 # ===== #
59 print('Computation time: %.2f'%(time.time()-time_start))

```

Fold done
 Fold done
 Fold done
 Fold done
 Fold done
 Done cross validation for norm l1. Error is 0.70
 Fold done
 Fold done
 Fold done
 Fold done
 Fold done
 Done cross validation for norm l2. Error is 0.73
 Fold done
 Fold done
 Fold done
 Fold done
 Done cross validation for norm inf. Error is 0.83
 Computation time: 711.44



In [64]: 1 normError

Out[64]: [0.7036, 0.7295999999999999, 0.8318]

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) It is 0.7036

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 [14]: 1 error = 1
          2
          3 # ===== #
          4 # YOUR CODE HERE:
          5 #   Evaluate the testing error of the k-nearest neighbors classifier
          6 #   for your optimal hyperparameters found by 5-fold cross-validation.
          7 # ===== #
          8 bestK=20
          9 bestKnn=KNN()
         10 bestKnn.train(X=X_train, y=y_train)
         11 distances = bestKnn.compute_distances(X=X_test,norm=L1_norm)
         12 predictions = bestKnn.predict_labels(distances,bestK)
         13 error = sum(predictions!=y_test)/len(y_test)
         14
         15 # ===== #
         16 # END YOUR CODE HERE
         17 # ===== #
         18
         19 print('Error rate achieved: {}'.format(error))
```

Error rate achieved: 0.72

Question:

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

Answer:

The naive model with $k=1$ had a test error of 0.726. By implementing cross validation and choosing the best k and norm values, we managed to reduce the error to 0.72. Therefore, we improved the performance by 0.06.