

knn_nosol

January 29, 2023

0.1 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.

0.2 Import the appropriate libraries

```
[61]: 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:

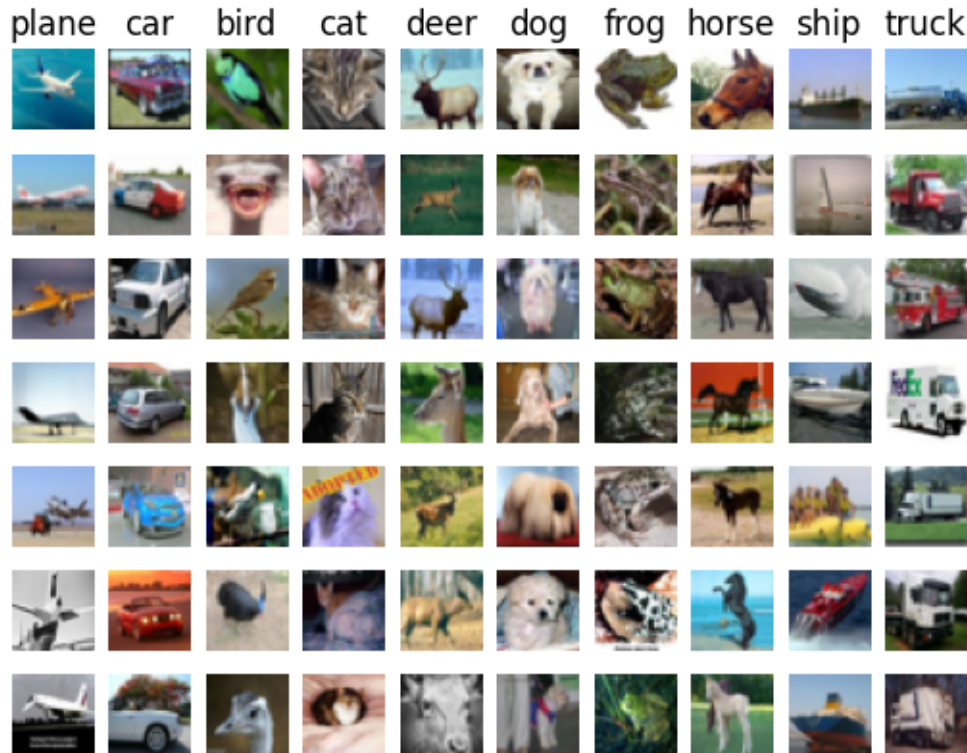
```
%reload_ext autoreload
```

```
[62]: # Set the path to the CIFAR-10 data
cifar10_dir = "/Users/mylesthemonster/Documents/ece_c247/hw2/hw2_code/
cifar-10-batches-py"
X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

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

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

```
[63]: # 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()
```



```
[64]: # 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)

1 K-nearest neighbors

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

```
[65]: # Import the KNN class
      from nn1 import KNN
```

```
[66]: # 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)
```

1.1 Questions

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

1.2 Answers

- (1) The inside the `knn.train()` function looks like:

```
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
```

All this being done is that the training data is being stored in the class.

- (2) The pros of this training step are that it is simple. A con of this training step is that it will take up memory.

1.3 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.

```
[67]: # Implement the function compute_distances() in the KNN class.
      # Do not worry about the input 'norm' for now; use the default definition of 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: 14.396498918533325
 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

1.3.1 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.

```
[68]: # 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: 0.06827902793884277
 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.

1.3.2 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.

```
[69]: # 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.
# ===== #
yPredicted = knn.predict_labels(dists_L2_vectorized, 1)
error = np.count_nonzero(y_test - yPredicted) / float(len(y_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%.

2 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.

2.0.1 Create training and validation folds

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

```

[70]: # 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]
# ===== #

# Calculate the size of each fold
n = X_train.shape[0]

```

```

# Divide the number of examples by the number of folds
fold_size = int(n / num_folds)

# Iterate over the number of folds
for i in range(num_folds):

    # Append each fold of the training data to the X_train_folds list
    X_train_folds.append(X_train[i * fold_size : i * fold_size + fold_size])

    # Append each fold of the corresponding labels to the y_train_folds list
    y_train_folds.append(y_train[i * fold_size : i * fold_size + fold_size])

print("==> y_train.shape: ", y_train.shape)
print("==> X_train.shape: ", X_train.shape)
print("==> y_train_folds[0].shape: ", y_train_folds[0].shape)
print("==> X_train_folds[0].shape: ", X_train_folds[0].shape)
print("==> Labels in each fold: ", y_train_folds[0].shape[0])
print("==> Training data entries in each fold: ", X_train_folds[0].shape[0])

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

```

```

==> y_train.shape: (5000,)
==> X_train.shape: (5000, 3072)
==> y_train_folds[0].shape: (1000,)
==> X_train_folds[0].shape: (1000, 3072)
==> Labels in each fold: 1000
==> Training data entries in each fold: 1000

```

2.0.2 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.

```

[75]: 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.
# ===== #

```

```

# List to store the average cross-validation error for each k
avr_cross_val_err = []
entries_per_fold = y_train_folds[0].shape[0]

# Loop through each k
for k in ks:
    total_error = 0

    # Loop through each fold
    for i in range(num_folds):
        # Declare an instance of the knn class.
        knn = KNN()

        # Create the training and testing sets for the current fold
        X_test_fold = X_train_folds[i]
        y_test_fold = y_train_folds[i]

        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 :
↪])

        # Train the model on the current training set
        knn.train(X=X_train_fold, y=y_train_fold)

        # Compute the L2 distances and predict the labels using the current
↪value of k
        dists_fold = knn.compute_L2_distances_vectorized(X_test_fold)
        y_est_fold = knn.predict_labels(dists_fold, k)

        # Calculate the number of correct predictions
        total_correct = np.sum(y_test_fold == y_est_fold)

        # Calculate the error for the current fold
        error = (entries_per_fold - total_correct) / entries_per_fold

        # Add the error for the current fold to the total error
        total_error += error

    # Append the average error for the current value of k to the list of errors
    avr_cross_val_err.append(total_error / num_folds)

index_min_error = np.argmin(avr_cross_val_err)
print(f"The optimal k is k = {ks[index_min_error]}, with a cross-validation
↪error of {avr_cross_val_err[index_min_error]}")

```



```

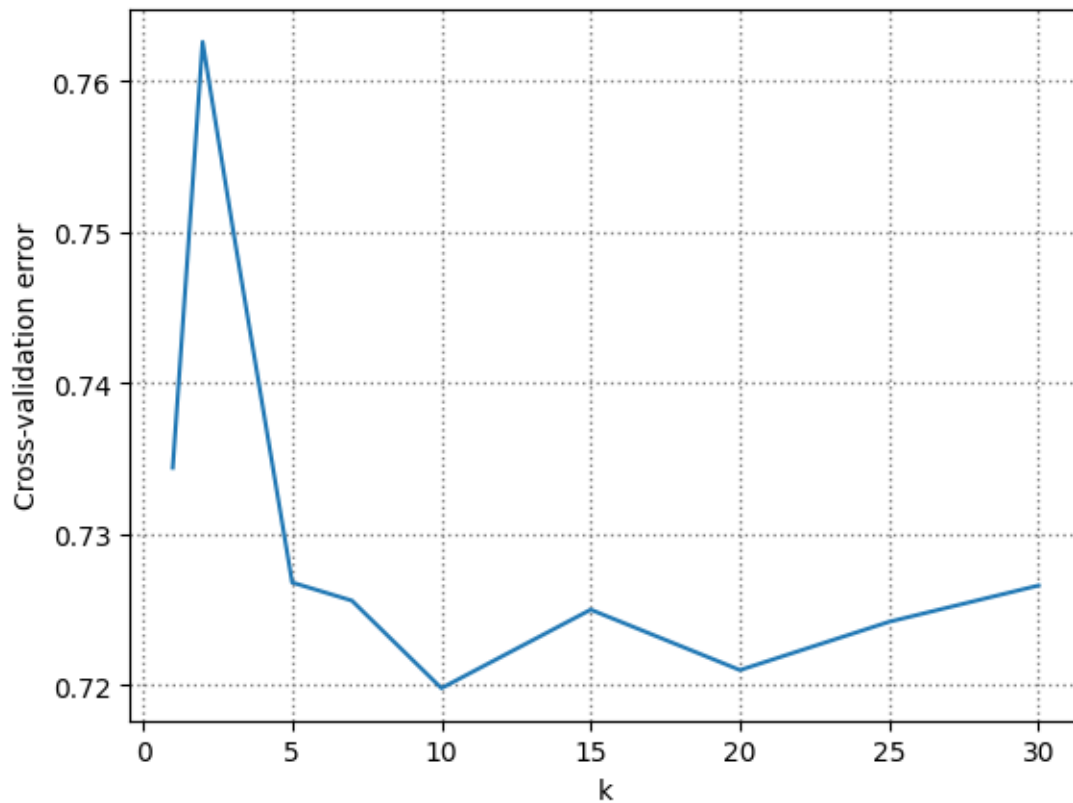
# Plot k vs. Cross-validation Error
plt.plot(ks, avr_cross_val_err)
plt.xlabel("k")
plt.ylabel("Cross-validation error")
plt.grid(color="grey", linestyle=":", linewidth=1)
plt.show()

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

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

```

The optimal k is $k = 10$, with a cross-validation error of 0.7198



Computation time: 14.99

2.1 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 ?

2.2 Answers:

- (1) The best value of k amongst the tested k 's is $k = 10$
- (2) The cross-validation error for $k = 10$ is 0.7198

2.2.1 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.

```
[72]: 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.
# ===== #

# Vectorized form of the L1- and Linf- norms
Vec_L1_norm = lambda x: np.sum(np.abs(x))
Vec_L2_norm = lambda x: np.sqrt(np.sum(x**2))
Vec_Linf_norm = lambda x: np.max(np.abs(x))
vec_norms = [Vec_L1_norm, Vec_L2_norm, Vec_Linf_norm]
vec_norms_names = ["L1", "L2", "Linf"]

# List to store the average cross-validation error for each norm
avr_cross_val_err = []
entries_per_fold = y_train_folds[0].shape[0]

# Best k from the previous part
k = 10

# Iterate over each norm
for l in vec_norms:
    total_error = 0
```

```

# Iterate over each fold
for i in range(num_folds):
    # Initialize KNN classifier
    knn = KNN()

    # Create the training and testing sets for the current fold
    X_test_fold = X_train_folds[i]
    y_test_fold = y_train_folds[i]

    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 :
↪])

    # Train the model on the current training set
    knn.train(X=X_train_fold, y=y_train_fold)

    # Compute the distances between the test data and the train data using
↪the current norm
    dists_fold = knn.compute_distances(X_test_fold, 1)

    # Predict the labels for the test data using the current norm
    y_est_fold = knn.predict_labels(dists_fold, k)

    # Calculate the error for the current fold
    y_diff_fold = y_test_fold - y_est_fold

    # Calculate the number of correct predictions
    total_correct = np.sum(y_test_fold == y_est_fold)

    # Calculate the error for the current fold
    error = (entries_per_fold - total_correct) / entries_per_fold

    # Add the error for the current fold to the total error
    total_error += error

    # Append the average error for the current value of k to the list of errors
    avr_cross_val_err.append(total_error / num_folds)

# Print the errors for each norm
for j in np.arange(len(avr_cross_val_err)):
    print(
        f"For the {vec_norms_names[j]} vectorized norm , the cross-validation
↪error is {avr_cross_val_err[j]}"
    )

# Plot the error vs the norm used

```

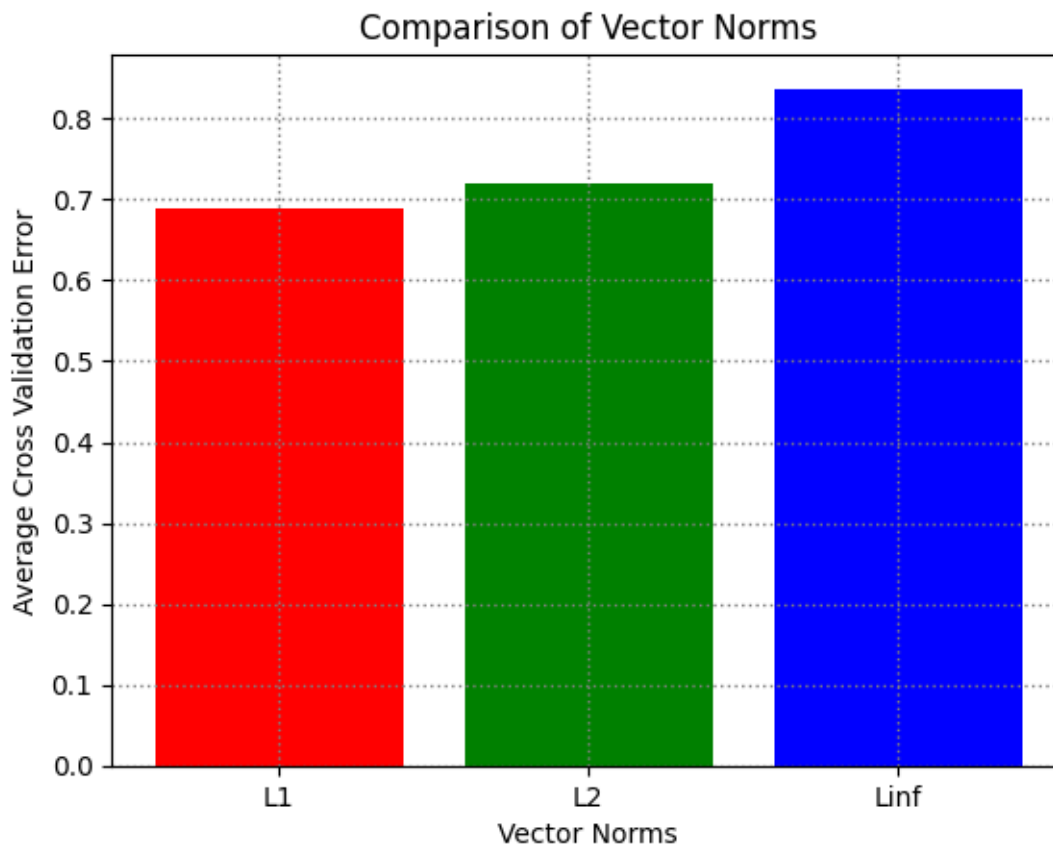
```

plt.figure()
plt.bar(vec_norms_names, avr_cross_val_err, color=['red', 'green', 'blue'])
plt.xlabel("Vector Norms")
plt.ylabel("Average Cross Validation Error")
plt.title("Comparison of Vector Norms")
plt.grid(color="grey", linestyle=":", linewidth=1)
plt.show()

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

```

For the L1 vectorized norm , the cross-validation error is 0.6886000000000001
For the L2 vectorized norm , the cross-validation error is 0.7198
For the Linf vectorized norm , the cross-validation error is 0.8370000000000001



Computation time: 311.81

2.3 Questions:

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

2.4 Answers:

- (1) The $L1$ norm has the best cross-validation error
- (2) the cross-validation error for the $L1$ norm and $k = 10$ is 0.6886000000000001

3 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.

```
[73]: 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_test, Vec_L1_norm)
y_pred = knn.predict_labels(dists_L1, k)
error = np.count_nonzero(y_test - y_pred) / float(len(y_test))

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

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

Error rate achieved: 0.722

3.1 Question:

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

3.2 Answer:

My error from by cross-validation with $k = 1$ and using the $L2$ -norm was 0.726 and my error from cross-validation with the optimal $k = 10$ and $L1$ -norm is 0.722. This means that my error improved by 0.004.