

knn

January 27, 2021

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.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (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.

0.2 Import the appropriate libraries

```
[8]: 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
    ↪ 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
```

```
[11]: # Set the path to the CIFAR-10 data
cifar10_dir = 'C:/Users/Ashwin/Desktop/UCLA/current classes/ece 247/hw2/
    ↪ hw2_code/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)

```

```

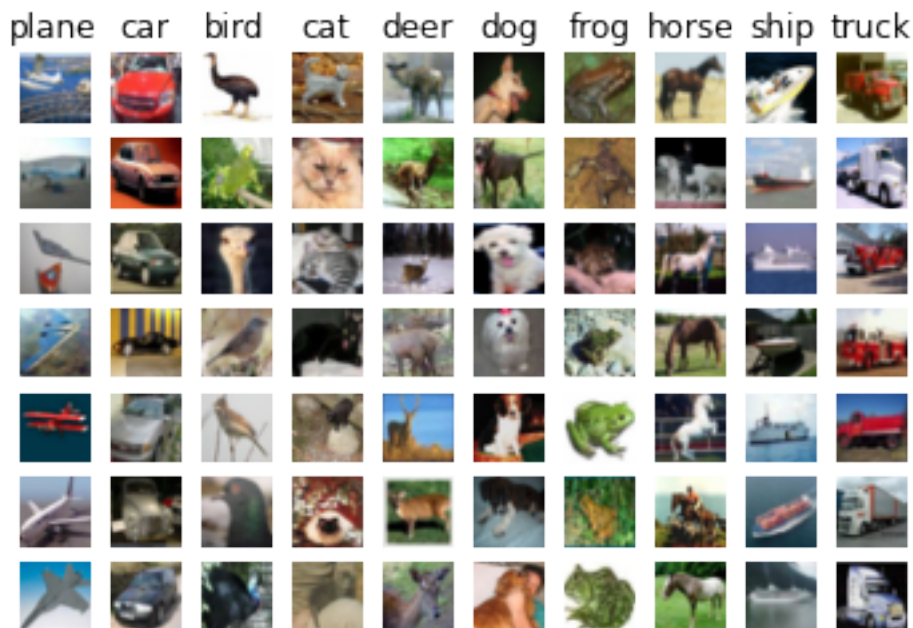
Training data shape: (50000, 32, 32, 3)
Training labels shape: (50000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)

```

```

[12]: # 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()

```



```
[13]: # 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.

```
[14]: # Import the KNN class

from nndl import KNN
```

```
[15]: # 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) In `knn.train()`, we are just declaring and initializing the instance variables `X_train` and `y_train` of the `knn` object and caching them in memory.

- (2) The pros is that it's fast, $O(1)$, and very simple. The biggest con is that it takes a lot of memory to store all the training and test data in the object.

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.

```
[47]: # 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: 105.75354480743408

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.

```
[30]: # 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: 0.6139991283416748

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.

```
[61]: # 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

# ===== #
# 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, 1)
num_incorrect = np.sum(y_test != predicted_labels)
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%.

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.

```
[66]: # 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 = num_training // num_folds
for i in range(num_folds):
    start, end = fold_size*i, fold_size*(i+1)
    X_train_folds.append(X_train[start:end])
    y_train_folds.append(y_train[start:end])

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

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.

```
[103]: time_start =time.time()

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

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

def calculate_error_per_k (k):
    errors = []

    for i, xy in enumerate(zip(X_train_folds, y_train_folds)):
```

```

xtest, ytest = xy[0], xy[1]

#training data for this fold = all training data - this fold
start, end = fold_size*i, fold_size*(i+1)

#kept getting errors when trying to use X_train and splicing, so I'll
→ just do it manually
xtrain_list = []
ytrain_list = []

for j in range(num_folds):
    if i == j:
        continue
    xtrain_list.extend(X_train_folds[j])
    ytrain_list.extend(y_train_folds[j])

#convert to numpy array since lists are not compatible
xtrain, ytrain = np.array(xtrain_list), np.array(ytrain_list)

#train model on the other k-1 folds
model = KNN()
model.train(X=xtrain, y=ytrain)

#test model on this specific fold
predicted_labels = model.predict_labels(model.
→ compute_L2_distances_vectorized(xtest), k)
num_incorrect = np.sum(ytest != predicted_labels)
error = num_incorrect / len(xtest)
errors.append(error)
return np.mean(errors)

error_per_ks = [calculate_error_per_k(x) for x in ks]

for x,y in zip(ks,error_per_ks):
    print("k:", x, "error:", y)

plt.scatter(ks, error_per_ks)
plt.show()
# ===== #
# END YOUR CODE HERE
# ===== #

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

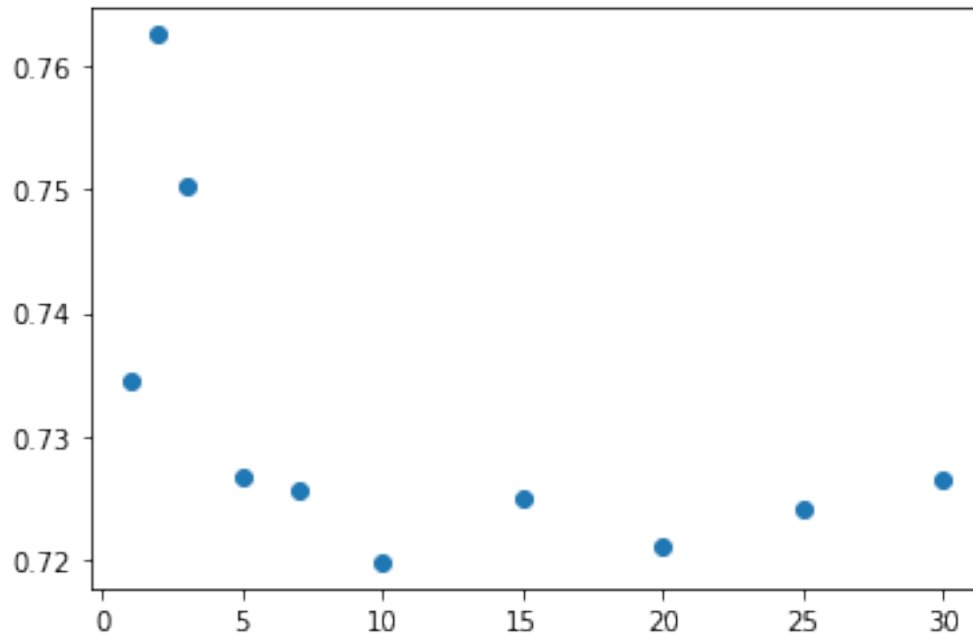
```

```

k: 1 error: 0.7344
k: 2 error: 0.76260000000000002
k: 3 error: 0.75040000000000001

```

k: 5 error: 0.7267999999999999
k: 7 error: 0.7256
k: 10 error: 0.7198
k: 15 error: 0.725
k: 20 error: 0.721
k: 25 error: 0.7242
k: 30 error: 0.7266



Computation time: 93.00

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) Looking at the results, it looks like $k = 10$ is the best among the tested k values.
- (2) The cross-validation error 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.


```

[104]: 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.
# ===== #

def calculate_error_per_k_norm (k, norm):
    errors = []

    for i, xy in enumerate(zip(X_train_folds, y_train_folds)):

        xtest, ytest = xy[0], xy[1]

        #training data for this fold = all training data - this fold
        start, end = fold_size*i, fold_size*(i+1)

        #kept getting errors when trying to use X_train and splicing, so I'll
→ just do it manually
        xtrain_list = []
        ytrain_list = []

        for j in range(num_folds):
            if i == j:
                continue
            xtrain_list.extend(X_train_folds[j])
            ytrain_list.extend(y_train_folds[j])

        #convert to numpy array since lists are not compatible
        xtrain, ytrain = np.array(xtrain_list), np.array(ytrain_list)

        #train model on the other k-1 folds
        model = KNN()
        model.train(X=xtrain, y=ytrain)

```

```

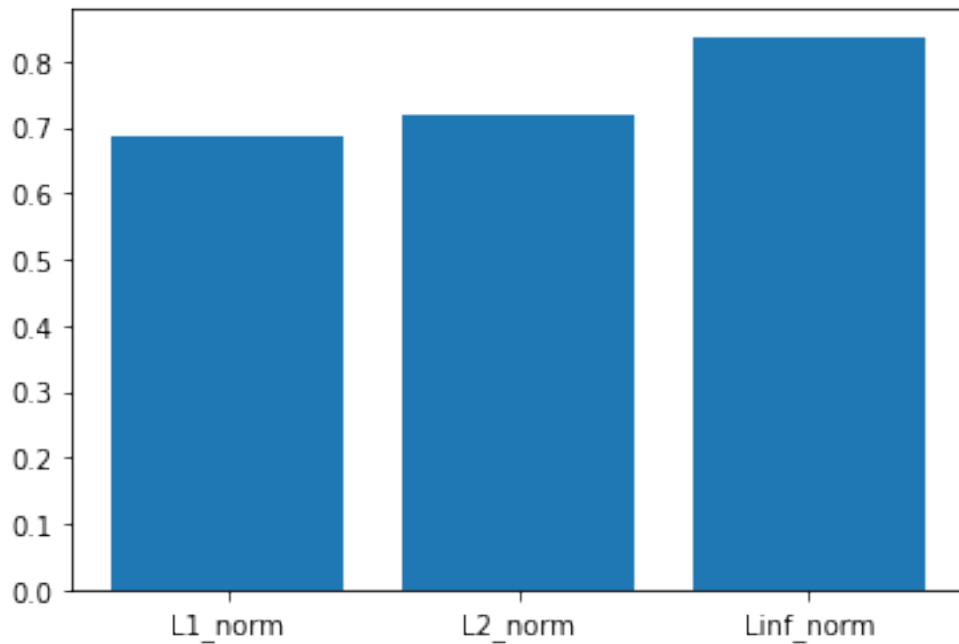
    #test model on this specific fold
    predicted_labels = model.predict_labels(model.compute_distances(xtest,
↪norm), k)
    num_incorrect = np.sum(ytest != predicted_labels)
    error = num_incorrect / len(xtest)
    errors.append(error)
    return np.mean(errors)

errors_per_norm = [calculate_error_per_k_norm(10,norm) for norm in norms]

plt.bar(["L1_norm", "L2_norm", "Linf_norm"], errors_per_norm)
plt.show()

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

```



Computation time: 1941.78

[105]: `print(errors_per_norm)`

[0.6886000000000001, 0.7198, 0.8370000000000001]

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) L1 norm – hopefully this is correct b/c it took me 30 min to run oops.
- (2) 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.

```
[107]: # ===== #
# YOUR CODE HERE:
# Evaluate the testing error of the k-nearest neighbors classifier
# for your optimal hyperparameters found by 5-fold cross-validation.
# ===== #

dist_l1 = knn.compute_distances(X=X_test, norm = L1_norm)
predicted_l1 = knn.predict_labels(dist_l1, 10)
error = np.sum(y_test != predicted_l1) / num_test

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

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

#difference between old error and new error
print("Improvement:", 0.726 - error)
```

```
Error rate achieved: 0.722
Improvement: 0.0040000000000000036
```

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:

The error rate improved slightly by about 0.004.

3.3 KNN Notebook Code Below:

```
[ ]: import numpy as np
import pdb
from scipy.stats import mode

"""
This code was based off of code from cs231n at Stanford University, and
↳ modified for ECE C147/C247 at UCLA.
"""

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])
        pass

        #
        ===== #
        # 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.
# ===== #
#let X be a N x 1 column vector, we want to expand it to a N x M matrix
#let X_train be a 1 x M row vector, we want to expand it to a N x M matrix

#optimized pairwise distances from here: https://www.pythonlikeyoumeanit.com/Module3\_IntroducingNumpy/Broadcasting.html#Pairwise-Distances-Using-Broadcasting-%28Unoptimized%29

x_squared = np.sum(X**2, axis=1)[:, np.newaxis]
y_squared = np.sum(self.X_train**2, axis=1)
two_x_y = 2*np.matmul(X, self.X_train.T)

dists = np.sqrt(x_squared + y_squared - two_x_y)

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

k_nearest = dists[i].argsort()[:k] #finds the indices

#convert indices -> labels
closest_y = self.y_train[k_nearest]

#find mode, aka the label that appears the most; pick the smaller label
→ to break ties
y_pred[i] = mode(closest_y)[0][0]

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

return y_pred

```