

# KNN Workbook for CS145 Homework 3

---

**\*\*PRINT YOUR NAME AND UID HERE!\*\***

NAME: [PUGALIA, GARVIT RAJKUMAR] UID: [504628127]

---

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 [148]: import numpy as np # for doing most of our calculations
import matplotlib.pyplot as plt # for plotting
from cs145.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-python
%load_ext autoreload
%autoreload 2
```

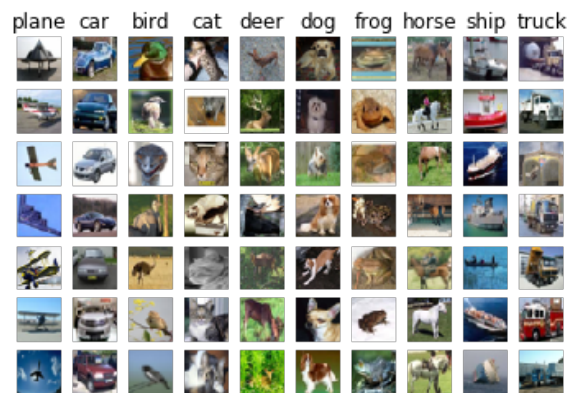
The autoreload extension is already loaded. To reload it, use:  
%reload\_ext autoreload

```
In [149]: # Set the path to the CIFAR-10 data
cifar10_dir = './cs145/datasets/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,)
```

```
In [150]: # 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 [151]: # 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 [152]: # Import the KNN class
from lib import KNN
```

```
In [153]: # 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 of KNN?

## Answers

- (1) In `knn.train()`, we store the entire dataset within the KNN classifier. This includes all features and respective labels for every data point. The classifier will use this stored information while calculating distances in the prediction phase.
- (2) The major pro is the lack of complex computation in the training phase. The `train()` function is simple, easy to implement, and fast under most circumstances. The con is the amount of memory consumed to store the entire dataset. The dataset is required to be held until prediction, making the classifier memory-inefficient.
- 

## 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 [154]: # Implement the function compute_distances() in the KNN class.
# Do not worry about the input 'norm' for now; use the default definition o
f 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: 45.794371128082275
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. Normally it may takes 20-40 seconds.

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 [155]: # Implement the function compute_L2_distances_vectorized() in the KNN class
          .
          # In this function, you ought to achieve the same L2 distance but WITHOUT a
          ny 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: 0.8920366764068604
          Difference in L2 distances between your KNN implementations (should be 0):
          1.4651847440245846e-10
```

## Speedup

Depending on your computer speed, you should see a 20-100x speed up from vectorization and no difference in L2 distances between two implementations. (There is a minute difference in the implementations, which is within acceptable range)

On our computer, the vectorized form took 0.20 seconds while the naive implementation took 26.88 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 [156]: # 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)
diff = y_test - y_pred
num_incorrect = np.sum(diff != 0)
error = num_incorrect / y_test.shape[0]

# ===== #
# 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.

## 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 [157]: # 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]
# ===== #
X_train_folds = np.array_split(X_train, num_folds)
y_train_folds = np.array_split(y_train, num_folds)

pass

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

X_train_folds = np.asarray(X_train_folds)
y_train_folds = np.asarray(y_train_folds)
```

## 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 [158]: time_start =time.time()

ks = [1, 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.
# ===== #
average_cross_validation_errors = []
for k_ in ks:
    error_for_run = []
    # ensure that the two datasets are shuffled in unison
    random_shuffle_state = np.random.get_state()
    np.random.shuffle(X_train_folds)
    np.random.set_state(random_shuffle_state)
    np.random.shuffle(y_train_folds)
    for i in range(num_folds):
        # separate training and testing data
        permut = [x for x in range(num_folds) if x != i]

        X_train_data = np.concatenate(X_train_folds[permut])
        y_train_data = np.concatenate(y_train_folds[permut])

        # generate a new classifier and train
        knn1 = KNN()
        knn1.train(X=X_train_data, y=y_train_data)
        dists_L2_vectorized_t = knn1.compute_L2_distances_vectorized(X=X_train_data, y=y_train_folds[i])
        y_pred = knn1.predict_labels(dists_L2_vectorized_t, k=k_)

        # calculate error for run
        error_for_run.append(float(np.sum(y_pred != y_train_folds[i])) / float(y_train_folds[i].shape[0]))
    average_cross_validation_errors.append(float(np.sum(error_for_run)) / float(num_folds))

results = np.zeros(len(ks))
results = average_cross_validation_errors

ks_min = ks[np.argsort(results)[0]]
results_min = min(results)

print('Set k = {0} and get minimum error as {1}'.format(ks_min,results_min))

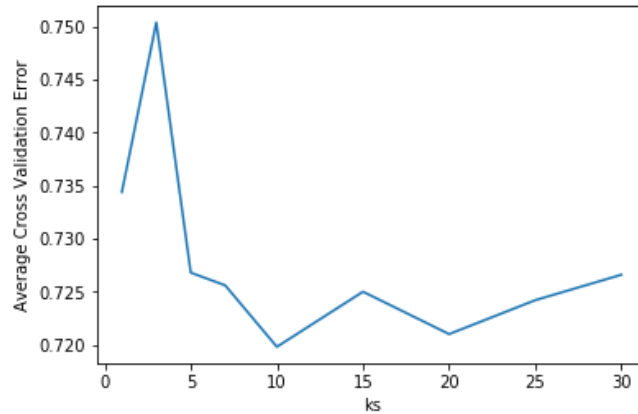
# Plot of k vs. cross-validation error
plt.plot(ks,results)
plt.xlabel('ks')
plt.ylabel('Average Cross Validation Error')
plt.show()

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

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

```

Set  $k = 10$  and get minimum error as 0.7198



Computation time: 100.79

### 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  $k$  (from the above experiment) is 10. (2) The cross-validation error for  $k = 10$  is 0.7198
- 

## 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 [159]: 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.
# ===== #

dists_L2_vectorized = knn.compute_L2_distances_vectorized(X=X_test)
y_pred = knn.predict_labels(dists_L2_vectorized, k=10)
diff = y_test - y_pred
num_incorrect = np.sum(diff != 0)
error = num_incorrect / y_test.shape[0]

# ===== #
# 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:

The error rate achieved with  $k = 10$  is 0.718, whereas the original naive implementation produced an error rate of 0.726. Different values of  $k$  had little impact on the cross-validation error with a 1.1% relative improvement with the best value of  $k = 10$ . While the actual code can be incorrect, I believe that the dataset might not be well-suited to a  $k$ -neighbors classifier (therefore, values of  $k$  have little impact). A change in norm or a change in approach is required to break the 27-30% accuracy barrier.

---

## The End of KNN Workbook

Please export this workbook as PDF file (see instructions) after completion.