

# knn

February 14, 2019

## 1 k-Nearest Neighbor (kNN) exercise

*Complete and hand in this completed worksheet (including its outputs and any supporting code outside of the worksheet) with your assignment submission.*

The kNN classifier consists of two stages:

- During training, the classifier takes the training data and simply remembers it
- During testing, kNN classifies every test image by comparing to all training images and transferring the labels of the k most similar training examples
- The value of k is cross-validated

In this exercise you will implement these steps and understand the basic Image Classification pipeline, cross-validation, and gain proficiency in writing efficient, vectorized code.

In [1]: *# Run some setup code for this notebook.*

```
from __future__ import print_function
import random
import numpy as np
from data_utils import load_CIFAR10
import matplotlib.pyplot as plt

# This is a bit of magic to make matplotlib figures appear inline in the notebook
# rather than in a new window.
%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# Some more magic so that the notebook will reload external python modules;
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2
```

```
In [2]: def rel_error(out, correct_out):
        return np.sum(abs(out - correct_out) / (abs(out) + abs(correct_out)))
```

```
In [3]: # Load the raw CIFAR-10 data.
cifar10_dir = '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 [4]: # 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 [5]: # Subsample the data for more efficient code execution in this exercise
num_training = 5000
mask = range(num_training)
```

```
X_train = X_train[mask]
y_train = y_train[mask]
```

```
num_test = 500
mask = range(num_test)
X_test = X_test[mask]
y_test = y_test[mask]
```

```
In [6]: # 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)
```

```
In [7]: from classifiers import KNearestNeighbor
```

```
# Create a kNN classifier instance.
# Remember that training a kNN classifier is a noop:
# the Classifier simply remembers the data and does no further processing
```

```
classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
```

We would now like to classify the test data with the kNN classifier. Recall that we can break down this process into two steps:

1. First we must compute the distances between all test examples and all train examples.
2. Given these distances, for each test example we find the  $k$  nearest examples and have them vote for the label

Lets begin with computing the distance matrix between all training and test examples. For example, if there are  $N_{tr}$  training examples and  $N_{te}$  test examples, this stage should result in a  $N_{te} \times N_{tr}$  matrix where each element  $(i,j)$  is the distance between the  $i$ -th test and  $j$ -th train example.

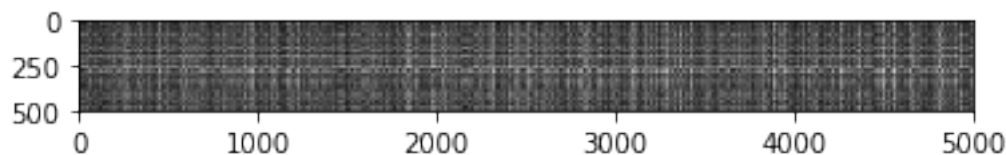
First, open `classifiers/k_nearest_neighbor.py` and implement the function `compute_distances_two_loops` that uses a (very inefficient) double loop over all pairs of (test, train) examples and computes the distance matrix one element at a time.

```
In [8]: # Open classifiers/k_nearest_neighbor.py and implement
        # compute_distances_two_loops.
```

```
# Test your implementation:
dists = classifier.compute_distances_two_loops(X_test)
print(dists.shape)
```

```
(500, 5000)
```

```
In [9]: # We can visualize the distance matrix: each row is a single test example and
        # its distances to training examples
plt.imshow(dists, interpolation='none')
plt.show()
```



**Inline Question #1:** Notice the structured patterns in the distance matrix, where some rows or columns are visible brighter. (Note that with the default color scheme black indicates low distances while white indicates high distances.)

- What in the data is the cause behind the distinctly bright rows?
- What causes the columns?

**Your Answer:** - This data point in testing data has high distances to the training data. - This data point in training data has high distances to the testing data.

```
In [10]: # Now implement the function predict_labels and run the code below:
# We use k = 1 (which is Nearest Neighbor).
y_test_pred = classifier.predict_labels(dists, k=1)

# Compute and print the fraction of correctly predicted examples
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 137 / 500 correct => accuracy: 0.274000

You should expect to see approximately 27% accuracy. Now let's try out a larger k, say k = 5:

```
In [11]: y_test_pred = classifier.predict_labels(dists, k=5)
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 139 / 500 correct => accuracy: 0.278000

You should expect to see a slightly better performance than with k = 1.

## 1.1 Frobenius Norm

To ensure that our vectorized implementation is correct, we make sure that it agrees with the naive implementation. There are many ways to decide whether two matrices are similar; one of the simplest is **the Frobenius norm**. \* Frobenius norm of  $m \times n$  matrix  $A$  is defined as the square root of the sum of the absolute squares of its elements,:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$$

```
In [12]: def Frobenius_norm(A):
    Fnorm = 0
    #####
    # TODO:
    # Implement a function to calculate Frobenius Norm of matrix A.
    # Hint: It is fine to use 2-nested for-loop. However, you can implement this
    # function with matrix calculation, which is much faster.
    # NOTE: numpy provides built-in function for Frobenius Norm, in this exercise,
    # you are required to implement this function.
    #####
    AA = np.dot(A, A.T)
    for i in range(AA.shape[0]):
        Fnorm += AA[i, i]
    Fnorm = np.sqrt(Fnorm)
```

```
#####
#                                     END OF YOUR CODE                                #
#####
return Fnorm
```

```
In [13]: # Check the accuracy of your implementation
A = np.random.rand(3,2)
print('The difference: ', rel_error(Frobenius_norm(A), np.linalg.norm(A)))
```

The difference: 0.0

```
In [14]: # Now lets speed up distance matrix computation by using partial vectorization
# with one loop. Implement the function compute_distances_one_loop and run the
# code below:
```

```
dists_one = classifier.compute_distances_one_loop(X_test)

difference = np.linalg.norm(dists - dists_one, ord='fro')
print('Difference was: %f' % (difference, ))
if difference < 0.001:
    print('Good! The distance matrices are the same')
else:
    print('Uh-oh! The distance matrices are different')
```

Difference was: 0.000000

Good! The distance matrices are the same

```
In [15]: # Now implement the fully vectorized version inside compute_distances_no_loops
# and run the code
```

```
dists_two = classifier.compute_distances_no_loops(X_test)
print('dists_two: ', dists_two)
print('dists: ', dists)
# check that the distance matrix agrees with the one we computed before:
difference = np.linalg.norm(dists - dists_two, ord='fro')
print('Difference was: %f' % (difference, ))
if difference < 0.001:
    print('Good! The distance matrices are the same')
else:
    print('Uh-oh! The distance matrices are different')
```

```
dists_two: [[3803.92350081 4210.59603857 5504.0544147 ... 4007.64756434
4203.28086142 4354.20256764]
[6336.83367306 5270.28006846 4040.63608854 ... 4829.15334194
4694.09767687 7768.33347636]
[5224.83913628 4250.64289255 3773.94581307 ... 3766.81549853
4464.99921613 6353.57190878]
...
[5366.93534524 5062.8772452 6361.85774755 ... 5126.56824786
```

```

4537.30613911 5920.94156364]
[3671.92919322 3858.60765044 4846.88157479 ... 3521.04515734
 3182.3673578 4448.65305458]
[6960.92443573 6083.71366848 6338.13442584 ... 6083.55504619
 4128.24744898 8041.05223214]]
dists: [[3803.92350081 4210.59603857 5504.0544147 ... 4007.64756434
 4203.28086142 4354.20256764]
[6336.83367306 5270.28006846 4040.63608854 ... 4829.15334194
 4694.09767687 7768.33347636]
[5224.83913628 4250.64289255 3773.94581307 ... 3766.81549853
 4464.99921613 6353.57190878]
...
[5366.93534524 5062.8772452 6361.85774755 ... 5126.56824786
 4537.30613911 5920.94156364]
[3671.92919322 3858.60765044 4846.88157479 ... 3521.04515734
 3182.3673578 4448.65305458]
[6960.92443573 6083.71366848 6338.13442584 ... 6083.55504619
 4128.24744898 8041.05223214]]
Difference was: 0.000000
Good! The distance matrices are the same

```

In [16]: *# Let's compare how fast the implementations are*

```

def time_function(f, *args):
    """
    Call a function f with args and return the time (in seconds) that it took to execute
    """
    import time
    tic = time.time()
    f(*args)
    toc = time.time()
    return toc - tic

two_loop_time = time_function(classifier.compute_distances_two_loops, X_test)
print('Two loop version took %f seconds' % two_loop_time)

one_loop_time = time_function(classifier.compute_distances_one_loop, X_test)
print('One loop version took %f seconds' % one_loop_time)

no_loop_time = time_function(classifier.compute_distances_no_loops, X_test)
print('No loop version took %f seconds' % no_loop_time)

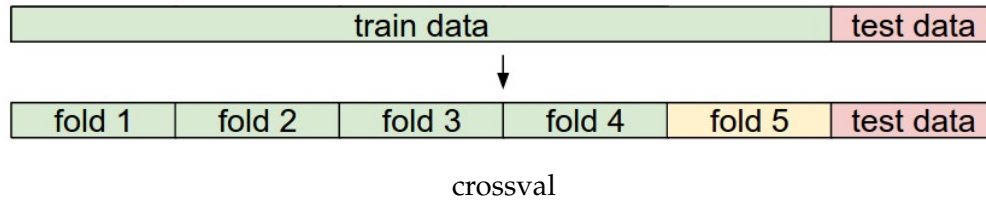
# you should see significantly faster performance with the fully vectorized implementation

```

```

Two loop version took 17.486984 seconds
One loop version took 19.122037 seconds
No loop version took 0.142319 seconds

```



### 1.1.1 Cross-validation

We have implemented the k-Nearest Neighbor classifier but we set the value  $k = 5$  arbitrarily. We will now determine the best value of this hyperparameter with cross-validation.

More detail explanation can be found [here](#).

```
In [17]: import copy
num_folds = 5
k_choices = [1, 3, 5, 8, 10, 12, 15, 20, 50, 100]

X_train_folds = []
y_train_folds = []
#####
# TODO:
# Split up the training data into folds. After splitting, X_train_folds and
# y_train_folds should each be lists of length num_folds, where
# y_train_folds[i] is the label vector for the points in X_train_folds[i].
# Hint: Look up the numpy array_split function.
#####
X_train_folds = np.array_split(X_train, num_folds)
y_train_folds = np.array_split(y_train, num_folds)
#####
#                                     END OF YOUR CODE                                     #
#####

# A dictionary holding the accuracies for different values of k that we find
# when running cross-validation. After running cross-validation,
# k_to_accuracies[k] should be a list of length num_folds giving the different
# accuracy values that we found when using that value of k.
k_to_accuracies = {}
for k_value in k_choices:
    k_to_accuracies[k_value] = []

    for i in range(5):
        Xk_val = X_train_folds[i]
        yk_val = y_train_folds[i]

        tempX_train = copy.deepcopy(X_train_folds)
        tempy_train = copy.deepcopy(y_train_folds)
        del tempX_train[i]
        del tempy_train[i]
```



```

Xk_train = np.concatenate(tempX_train)
yk_train = np.concatenate(tempy_train)

classifier_k = KNearestNeighbor()
classifier_k.train(Xk_train, yk_train)

dists = classifier_k.compute_distances_no_loops(Xk_val)
y_test_pred = classifier_k.predict_labels(dists, k=k_value)
num_correct = np.sum(y_test_pred == yk_val)
accuracy = float(num_correct) / len(Xk_val)
k_to_accuracies[k_value].append(accuracy)

#####
# TODO:
# Perform k-fold cross validation to find the best value of k. For each
# possible value of k, run the k-nearest-neighbor algorithm num_folds times,
# where in each case you use all but one of the folds as training data and the
# last fold as a validation set. Store the accuracies for all fold and all
# values of k in the k_to_accuracies dictionary.
#####

#####
#                                     END OF YOUR CODE
#####

# Print out the computed accuracies
for k in sorted(k_to_accuracies):
    for accuracy in k_to_accuracies[k]:
        print('k = %d, accuracy = %f' % (k, accuracy))

k = 1, accuracy = 0.263000
k = 1, accuracy = 0.257000
k = 1, accuracy = 0.264000
k = 1, accuracy = 0.278000
k = 1, accuracy = 0.266000
k = 3, accuracy = 0.241000
k = 3, accuracy = 0.249000
k = 3, accuracy = 0.243000
k = 3, accuracy = 0.273000
k = 3, accuracy = 0.264000
k = 5, accuracy = 0.256000
k = 5, accuracy = 0.271000
k = 5, accuracy = 0.280000
k = 5, accuracy = 0.289000
k = 5, accuracy = 0.278000
k = 8, accuracy = 0.263000
k = 8, accuracy = 0.287000
k = 8, accuracy = 0.276000

```

```

k = 8, accuracy = 0.288000
k = 8, accuracy = 0.270000
k = 10, accuracy = 0.266000
k = 10, accuracy = 0.296000
k = 10, accuracy = 0.279000
k = 10, accuracy = 0.283000
k = 10, accuracy = 0.283000
k = 12, accuracy = 0.261000
k = 12, accuracy = 0.294000
k = 12, accuracy = 0.280000
k = 12, accuracy = 0.283000
k = 12, accuracy = 0.280000
k = 15, accuracy = 0.253000
k = 15, accuracy = 0.290000
k = 15, accuracy = 0.279000
k = 15, accuracy = 0.280000
k = 15, accuracy = 0.275000
k = 20, accuracy = 0.270000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.280000
k = 20, accuracy = 0.284000
k = 50, accuracy = 0.271000
k = 50, accuracy = 0.288000
k = 50, accuracy = 0.278000
k = 50, accuracy = 0.269000
k = 50, accuracy = 0.266000
k = 100, accuracy = 0.256000
k = 100, accuracy = 0.270000
k = 100, accuracy = 0.263000
k = 100, accuracy = 0.256000
k = 100, accuracy = 0.263000

```

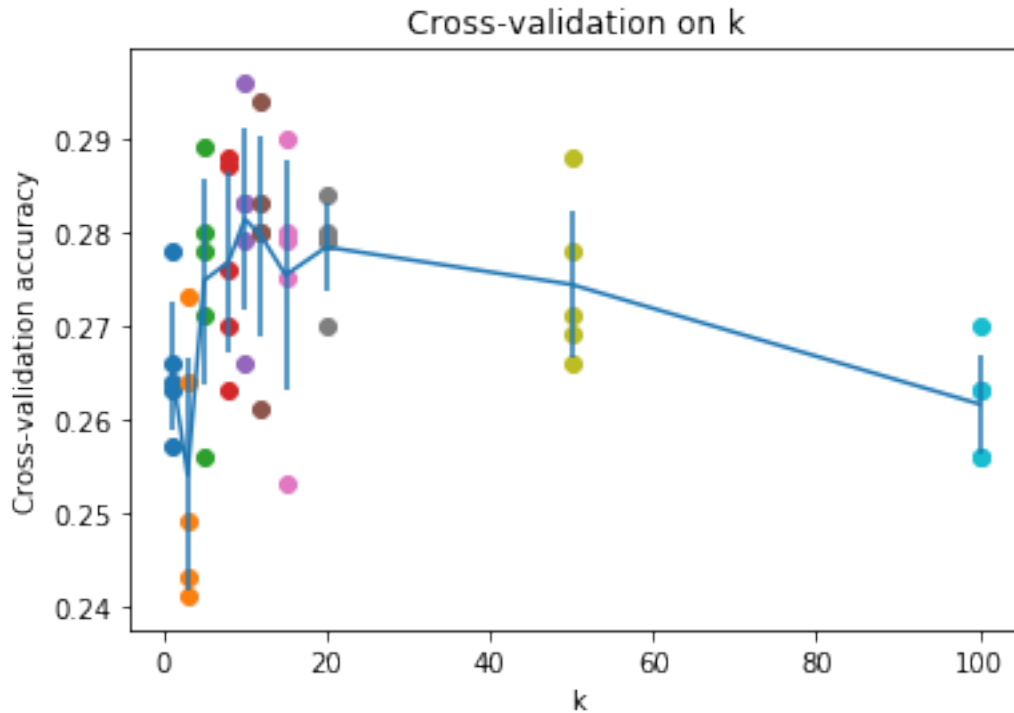
```
In [18]: # plot the raw observations
```

```

for k in k_choices:
    accuracies = k_to_accuracies[k]
    plt.scatter([k] * len(accuracies), accuracies)

# plot the trend line with error bars that correspond to standard deviation
accuracies_mean = np.array([np.mean(v) for k,v in sorted(k_to_accuracies.items())])
accuracies_std = np.array([np.std(v) for k,v in sorted(k_to_accuracies.items())])
plt.errorbar(k_choices, accuracies_mean, yerr=accuracies_std)
plt.title('Cross-validation on k')
plt.xlabel('k')
plt.ylabel('Cross-validation accuracy')
plt.show()

```



```
In [19]: # Based on the cross-validation results above, choose the best value for k,
# retrain the classifier using all the training data, and test it on the test
# data. You should be able to get above 10% accuracy on the test data.
best_k = 10

classifier = KNearestNeighbor()
classifier.train(X_train, y_train)
y_test_pred = classifier.predict(X_test, k=best_k, num_loops = 0)

# Compute and display the accuracy
num_correct = np.sum(y_test_pred == y_test)
accuracy = float(num_correct) / num_test
print ('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))

Got 141 / 500 correct => accuracy: 0.282000
```