

# 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.

**Acknowledgement: This exercise is adapted from [Stanford CS231n](#).**

```
In [52]: # 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

The autoreload extension is already loaded. To reload it, use:
%reload_ext autoreload
```

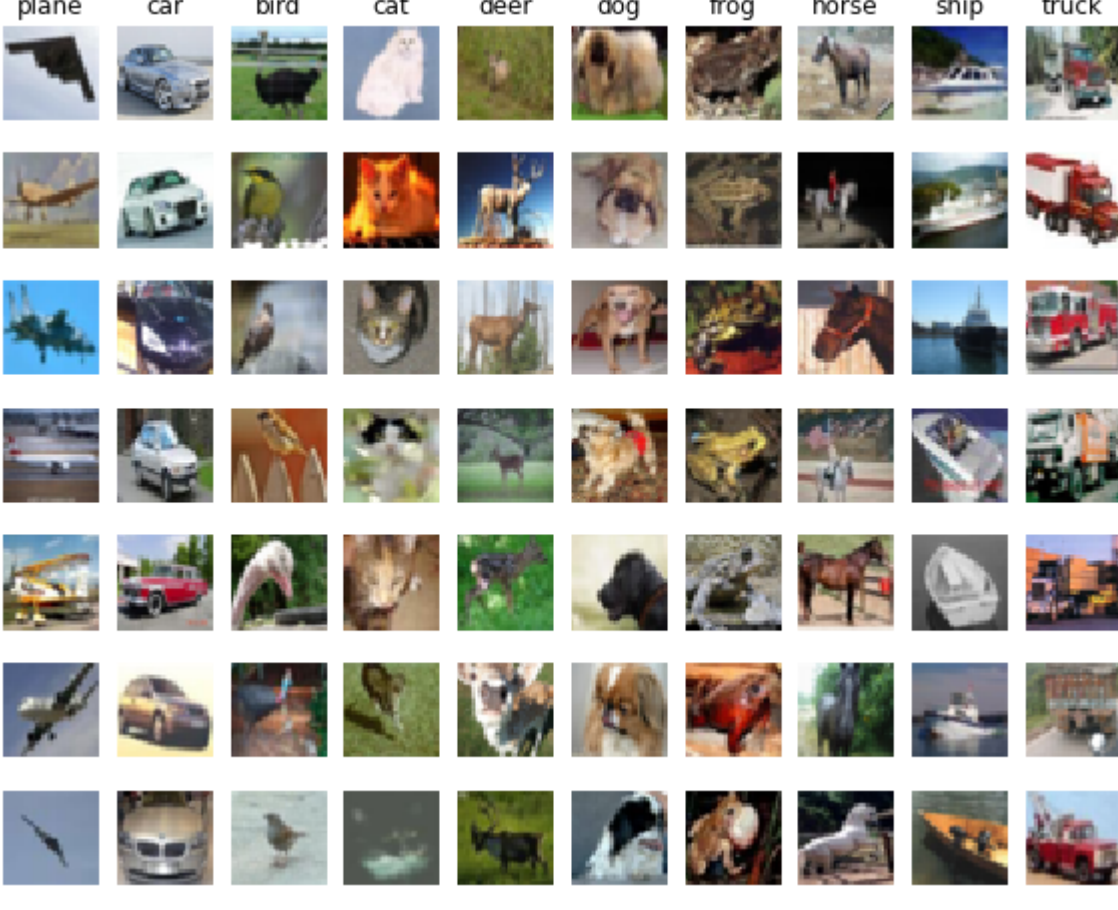
```
In [53]: def rel_error(out, correct_out):
        return np.sum(abs(out - correct_out) / (abs(out) + abs(correct_out)))
```

```
In [54]: # 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 [55]: # 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 [56]: # 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 [57]: # 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 [58]: 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 **Ntr** training examples and **Nte** test examples, this stage should result in a **Nte x Ntr** 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 [59]: # 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 [60]: # 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:** Test images which have a large distance between themselves and most training images cause distinctly bright rows. Each bright column represents a training example which has a large distance compared to all the test examples.

```
In [61]: # 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 lets try out a larger `k`, say `k = 5` :

```
In [62]: 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` .

## 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 [63]: def Frobenius_norm(A):
        Fnorm = None
        #####
        # 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.
        #####
        Fnorm = np.sqrt(np.sum(np.square(A)))
        #####
        # END OF YOUR CODE
        #####
        return Fnorm
```

```
In [64]: # 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 [65]: # 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 [66]: # 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 [67]: # 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 57.822969 seconds
One loop version took 85.534709 seconds
No loop version took 0.576198 seconds
```