

k-Nearest Neighbor (kNN) exercise

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]:

```
1 # Run some setup code for this notebook.
2
3 import random
4 import numpy as np
5 from cs175.data_utils import load_CIFAR10
6 import matplotlib.pyplot as plt
7
8 #from __future__ import print_function
9
10 # This is a bit of magic to make matplotlib figures appear inline in the notebook
11 # rather than in a new window.
12 %matplotlib inline
13 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
14 plt.rcParams['image.interpolation'] = 'nearest'
15 plt.rcParams['image.cmap'] = 'gray'
16
17 # Some more magic so that the notebook will reload external python modules;
18 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
19 %load_ext autoreload
20 %autoreload 2
```

In [2]:

```
1 # Load the raw CIFAR-10 data.
2 cifar10_dir = 'cs175/datasets/cifar-10-batches-py'
3 X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
4
5 # As a sanity check, we print out the size of the training and test data.
6 print('Training data shape: ', X_train.shape)
7 print('Training labels shape: ', y_train.shape)
8 print('Test data shape: ', X_test.shape)
9 print('Test labels shape: ', y_test.shape)
```

Training data shape: (50000, 32, 32, 3)

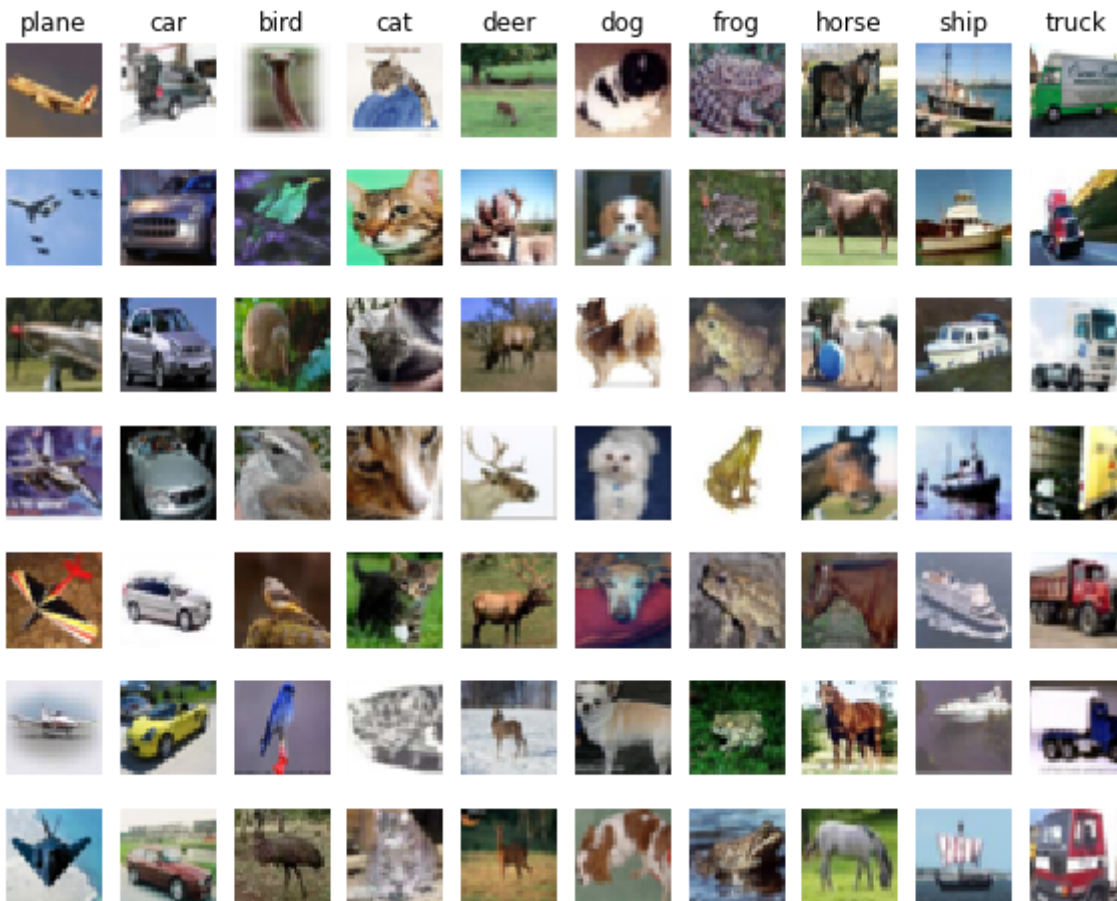
Training labels shape: (50000,)

Test data shape: (10000, 32, 32, 3)

Test labels shape: (10000,)

In [3]:

```
1 # Visualize some examples from the dataset.
2 # We show a few examples of training images from each class.
3 classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
4 num_classes = len(classes)
5 samples_per_class = 7
6 for y, cls in enumerate(classes):
7     idxs = np.flatnonzero(y_train == y)
8     idxs = np.random.choice(idxs, samples_per_class, replace=False)
9     for i, idx in enumerate(idxs):
10         plt_idx = i * num_classes + y + 1
11         plt.subplot(samples_per_class, num_classes, plt_idx)
12         plt.imshow(X_train[idx].astype('uint8'))
13         plt.axis('off')
14         if i == 0:
15             plt.title(cls)
16 plt.show()
```



In [4]:

```
1 # Subsample the data for more efficient code execution in this exercise
2 num_training = 5000
3 mask = list(range(num_training))
4 X_train = X_train[mask]
5 y_train = y_train[mask]
6
7 num_test = 500
8 mask = list(range(num_test))
9 X_test = X_test[mask]
10 y_test = y_test[mask]
```

In [5]:

```
1 # Reshape the image data into rows
2 X_train = np.reshape(X_train, (X_train.shape[0], -1))
3 X_test = np.reshape(X_test, (X_test.shape[0], -1))
4 print(X_train.shape, X_test.shape)
```

(5000, 3072) (500, 3072)

In [6]:

```
1 from cs175.classifiers import KNearestNeighbor
2
3 # Create a kNN classifier instance.
4 # Remember that training a kNN classifier is a noop:
5 # the Classifier simply remembers the data and does no further processing
6 classifier = KNearestNeighbor()
7 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 `cs175/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.

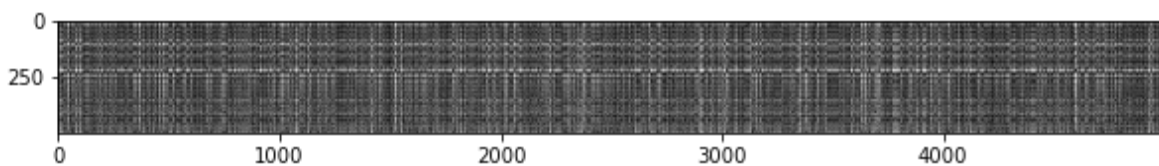
In [7]:

```
1 # Open cs175/classifiers/k_nearest_neighbor.py and implement
2 # compute_distances_two_loops.
3
4 # Test your implementation:
5 dists = classifier.compute_distances_two_loops(X_test)
6 print(dists.shape)
```

(500, 5000)

In [8]:

```
1 # We can visualize the distance matrix: each row is a single test example and
2 # its distances to training examples
3 plt.imshow(dists, interpolation='none')
4 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: *fill this in.* bright rows shows when some data in the test images has large distance(not similar) away from training images, therefore the large distance is shown by bright rows.

columns represents the training images has large distance(not similar) away from test images.

In [9]:

```
1 # Now implement the function predict_labels and run the code below:
2 # We use k = 1 (which is Nearest Neighbor).
3 y_test_pred = classifier.predict_labels(dists, k=1)
4
5 # Compute and print the fraction of correctly predicted examples
6 num_correct = np.sum(y_test_pred == y_test)
7 accuracy = float(num_correct) / num_test
8 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 [10]:

```
1 y_test_pred = classifier.predict_labels(dists, k=5)
2 num_correct = np.sum(y_test_pred == y_test)
3 accuracy = float(num_correct) / num_test
4 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` .

In [11]:

```
1 # Now lets speed up distance matrix computation by using partial vectorization
2 # with one loop. Implement the function compute_distances_one_loop and run the
3 # code below:
4 dists_one = classifier.compute_distances_one_loop(X_test)
5
6 # To ensure that our vectorized implementation is correct, we make sure that it
7 # agrees with the naive implementation. There are many ways to decide whether
8 # two matrices are similar; one of the simplest is the Frobenius norm. In case
9 # you haven't seen it before, the Frobenius norm of two matrices is the square
10 # root of the squared sum of differences of all elements; in other words, reshape
11 # the matrices into vectors and compute the Euclidean distance between them.
12 difference = np.linalg.norm(dists - dists_one, ord='fro')
13 print('Difference was: %f' % (difference, ))
14 if difference < 0.001:
15     print('Good! The distance matrices are the same')
16 else:
17     print('Uh-oh! The distance matrices are different')
```

Difference was: 0.000000

Good! The distance matrices are the same

In [12]:

```
1 # Now implement the fully vectorized version inside compute_distances_no_loops
2 # and run the code
3 dists_two = classifier.compute_distances_no_loops(X_test)
4
5 # check that the distance matrix agrees with the one we computed before:
6 difference = np.linalg.norm(dists - dists_two, ord='fro')
7 print('Difference was: %f' % (difference, ))
8 if difference < 0.001:
9     print('Good! The distance matrices are the same')
10 else:
11     print('Uh-oh! The distance matrices are different')
```

Difference was: 0.000000

Good! The distance matrices are the same

In [13]:

```
1 # Let's compare how fast the implementations are
2 def time_function(f, *args):
3     """
4     Call a function f with args and return the time (in seconds) that it took to execute.
5     """
6     import time
7     tic = time.time()
8     f(*args)
9     toc = time.time()
10    return toc - tic
11
12 two_loop_time = time_function(classifier.compute_distances_two_loops, X_test)
13 print('Two loop version took %f seconds' % two_loop_time)
14
15 one_loop_time = time_function(classifier.compute_distances_one_loop, X_test)
16 print('One loop version took %f seconds' % one_loop_time)
17
18 no_loop_time = time_function(classifier.compute_distances_no_loops, X_test)
19 print('No loop version took %f seconds' % no_loop_time)
20
21 # you should see significantly faster performance with the fully vectorized implementation
```

Two loop version took 24.351843 seconds

One loop version took 42.500599 seconds

No loop version took 0.159026 seconds

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.

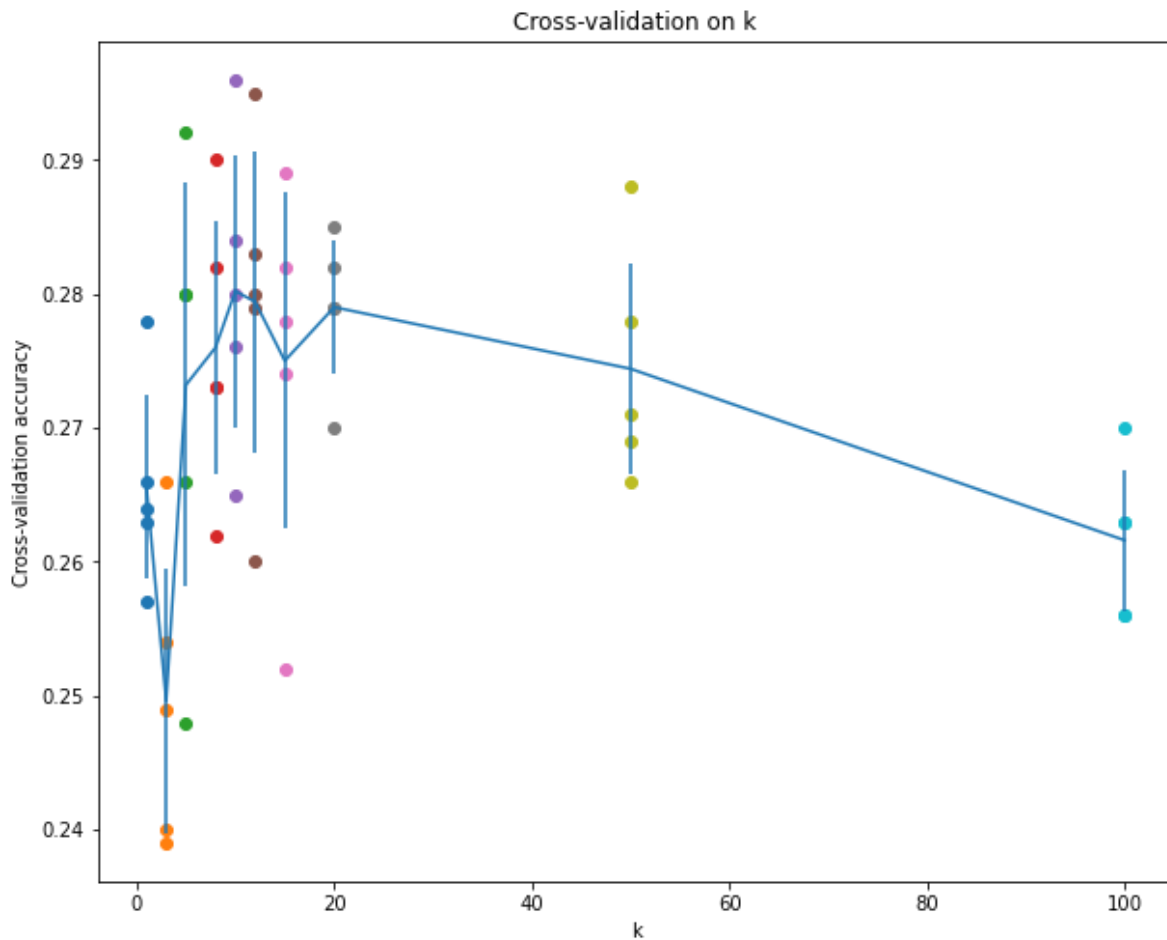
In [14]:

```
1 num_folds = 5
2 k_choices = [1, 3, 5, 8, 10, 12, 15, 20, 50, 100]
3
4 X_train_folds = []
5 y_train_folds = []
6 #####
7 # TODO: #
8 # Split up the training data into folds. After splitting, X_train_folds and #
9 # y_train_folds should each be lists of length num_folds, where #
10 # y_train_folds[i] is the label vector for the points in X_train_folds[i]. #
11 # Hint: Look up the numpy array_split function. #
12 #####
13 X_train_folds = np.array_split(X_train,num_folds)
14 y_train_folds = np.array_split(y_train,num_folds)
15 #####
16 #                               END OF YOUR CODE #
17 #####
18
19 # A dictionary holding the accuracies for different values of k that we find
20 # when running cross-validation. After running cross-validation,
21 # k_to_accuracies[k] should be a list of length num_folds giving the different
22 # accuracy values that we found when using that value of k.
23 k_to_accuracies = {}
24
25
26 num1 = X_train.shape[0]//num_folds
27 #####
28 # TODO: #
29 # Perform k-fold cross validation to find the best value of k. For each #
30 # possible value of k, run the k-nearest-neighbor algorithm num_folds times, #
31 # where in each case you use all but one of the folds as training data and the #
32 # last fold as a validation set. Store the accuracies for all fold and all #
33 # values of k in the k_to_accuracies dictionary. #
34 #####
35 for i in k_choices:
36     l1 = []
37     for j in range(num_folds):
38         Xtr = np.concatenate(X_train_folds[:j] + X_train_folds[j+1:])
39         Ytr = np.concatenate(y_train_folds[:j] + y_train_folds[j+1:])
40         Xva = X_train_folds[j]
41         Yva = y_train_folds[j]
42         classifier = KNearestNeighbor()
43         classifier.train(Xtr,Ytr)
44         Ypred = classifier.predict(Xva, k=i)
45         score= np.sum(Ypred == Yva)
46         accuracy = float(score)/num1
47
48         l1.append(accuracy)
49
50
51     k_to_accuracies[i] = l1
52
53 #####
54 #                               END OF YOUR CODE #
55 #####
56
57 # Print out the computed accuracies
58 for k in sorted(k_to_accuracies):
59     for accuracy in k_to_accuracies[k]:
```

```
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 = 0.239000
k = 3, accuracy = 0.249000
k = 3, accuracy = 0.240000
k = 3, accuracy = 0.266000
k = 3, accuracy = 0.254000
k = 5, accuracy = 0.248000
k = 5, accuracy = 0.266000
k = 5, accuracy = 0.280000
k = 5, accuracy = 0.292000
k = 5, accuracy = 0.280000
k = 8, accuracy = 0.262000
k = 8, accuracy = 0.282000
k = 8, accuracy = 0.273000
k = 8, accuracy = 0.290000
k = 8, accuracy = 0.273000
k = 10, accuracy = 0.265000
k = 10, accuracy = 0.296000
k = 10, accuracy = 0.276000
k = 10, accuracy = 0.284000
k = 10, accuracy = 0.280000
k = 12, accuracy = 0.260000
k = 12, accuracy = 0.295000
k = 12, accuracy = 0.279000
k = 12, accuracy = 0.283000
k = 12, accuracy = 0.280000
k = 15, accuracy = 0.252000
k = 15, accuracy = 0.289000
k = 15, accuracy = 0.278000
k = 15, accuracy = 0.282000
k = 15, accuracy = 0.274000
k = 20, accuracy = 0.270000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.279000
k = 20, accuracy = 0.282000
k = 20, accuracy = 0.285000
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 [15]:

```
1 # plot the raw observations
2 for k in k_choices:
3     accuracies = k_to_accuracies[k]
4     plt.scatter([k] * len(accuracies), accuracies)
5
6 # plot the trend line with error bars that correspond to standard deviation
7 accuracies_mean = np.array([np.mean(v) for k,v in sorted(k_to_accuracies.items())])
8 accuracies_std = np.array([np.std(v) for k,v in sorted(k_to_accuracies.items())])
9 plt.errorbar(k_choices, accuracies_mean, yerr=accuracies_std)
10 plt.title('Cross-validation on k')
11 plt.xlabel('k')
12 plt.ylabel('Cross-validation accuracy')
13 plt.show()
```



In [16]:

```
1 # Based on the cross-validation results above, choose the best value for k,
2 # retrain the classifier using all the training data, and test it on the test
3 # data. You should be able to get above 28% accuracy on the test data.
4 best_k = 10
5
6 classifier = KNearestNeighbor()
7 classifier.train(X_train, y_train)
8 y_test_pred = classifier.predict(X_test, k=best_k)
9
10 # Compute and display the accuracy
11 num_correct = np.sum(y_test_pred == y_test)
12 accuracy = float(num_correct) / num_test
13 print('Got %d / %d correct => accuracy: %f' % (num_correct, num_test, accuracy))
```

Got 141 / 500 correct => accuracy: 0.282000