

▼ EECS 498-007/598-005 Assignment 1-2: K-Nearest Neighbors (k-NN)

Before we start, please put your name and UMID in following format

: Firstname LASTNAME, #00000000 // e.g.) Justin JOHNSON, #12345678

Your Answer:

Artem KARPOV

K-Nearest Neighbors (k-NN)

In this notebook you will implement a K-Nearest Neighbors classifier on the [CIFAR-10 dataset](#).

Recall that the K-Nearest Neighbor classifier does the following:

- During training, the classifier simply memorizes the training data
- During testing, test images are compared to each training image; the predicted label is the majority vote among the K nearest training examples.

After implementing the K-Nearest Neighbor classifier, you will use *cross-validation* to find the best value of K.

The goals of this exercise are to go through a simple example of the data-driven image classification pipeline, and also to practice writing efficient, vectorized code in [PyTorch](#).

▼ Install starter code

We have implemented some utility functions for this exercise in the [cutils package](#). Run this cell to download and install it.

```
1 !pip install git+https://github.com/deepvision-class/starter-code
```

Looking in indexes: <https://pypi.org/simple>, <https://us-python.pkg.dev/colab-wheels/public/simple/>

Collecting git+<https://github.com/deepvision-class/starter-code>

Cloning <https://github.com/deepvision-class/starter-code> to /tmp/pip-req-build-j9e_202e

Running command git clone -q <https://github.com/deepvision-class/starter-code> /tmp/pip-req-build-j9e_202e

Requirement already satisfied: pydrive in /usr/local/lib/python3.7/dist-packages (from Colab-Utils==0.1.dev0) (1.3.1)

Requirement already satisfied: oauth2client>=4.0.0 in /usr/local/lib/python3.7/dist-packages (from pydrive->Colab-Utils==0.1.dev0) (4.1.3)

✓ 1 сек. выполнено в 08:13



```
Requirement already satisfied: pyyaml>=3.0 in /usr/local/lib/python3.7/dist-packages (from pydrive->Colab-Utills==0.1.dev0) (3.13)
Requirement already satisfied: httplib2<1dev,>=0.15.0 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive->Colab-
Requirement already satisfied: six<2dev,>=1.13.0 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive->Colab-Utills
Requirement already satisfied: google-auth-httplib2>=0.0.3 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive->C
Requirement already satisfied: google-auth<3dev,>=1.16.0 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive->Col
Requirement already satisfied: uritemplate<4dev,>=3.0.0 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive->Cola
Requirement already satisfied: google-api-core<3dev,>=1.21.0 in /usr/local/lib/python3.7/dist-packages (from google-api-python-client>=1.2->pydrive-
Requirement already satisfied: protobuf<4.0.0dev,>=3.12.0 in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0->google-api-
Requirement already satisfied: setuptools>=40.3.0 in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0->google-api-python-c
Requirement already satisfied: pytz in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0->google-api-python-client>=1.2->py
Requirement already satisfied: googleapis-common-protos<2.0dev,>=1.6.0 in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0
Requirement already satisfied: requests<3.0.0dev,>=2.18.0 in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0->google-api-
Requirement already satisfied: packaging>=14.3 in /usr/local/lib/python3.7/dist-packages (from google-api-core<3dev,>=1.21.0->google-api-python-clie
Requirement already satisfied: rsa<5,>=3.1.4 in /usr/local/lib/python3.7/dist-packages (from google-auth<3dev,>=1.16.0->google-api-python-client>=1.
Requirement already satisfied: cachetools<5.0,>=2.0.0 in /usr/local/lib/python3.7/dist-packages (from google-auth<3dev,>=1.16.0->google-api-python-c
Requirement already satisfied: pyasn1-modules>=0.2.1 in /usr/local/lib/python3.7/dist-packages (from google-auth<3dev,>=1.16.0->google-api-python-cl
Requirement already satisfied: pyasn1>=0.1.7 in /usr/local/lib/python3.7/dist-packages (from oauth2client>=4.0.0->pydrive->Colab-Utills==0.1.dev0) (0
Requirement already satisfied: pyparsing!=3.0.5,>=2.0.2 in /usr/local/lib/python3.7/dist-packages (from packaging>=14.3->google-api-core<3dev,>=1.21
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0dev,>=2.18.0->google-api-core<3dev,>=
Requirement already satisfied: urllib3!=1.25.0,!1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0dev,>=2.18.0->g
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0dev,>=2.18.0->google-api-core<3dev,>
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from requests<3.0.0dev,>=2.18.0->google-api-core<3dev,>=1.21.
```

▼ Setup code

Run some setup code for this notebook: Import some useful packages and increase the default figure size.

```
1 import couils
2 import torch
3 import torchvision
4 import matplotlib.pyplot as plt
5 import statistics
6
7 plt.rcParams['figure.figsize'] = (10.0, 8.0)
8 plt.rcParams['font.size'] = 16
```

▼ Load the CIFAR-10 dataset

The utility function `couils.data.cifar10()` returns the entire CIFAR-10 dataset as a set of four **Torch tensors**:

- `x_train` contains all training images (real numbers in the range $[0, 1]$)

- `y_train` contains all training labels (integers in the range `[0, 9]`)
- `x_test` contains all test images
- `y_test` contains all test labels

This function automatically downloads the CIFAR-10 dataset the first time you run it.

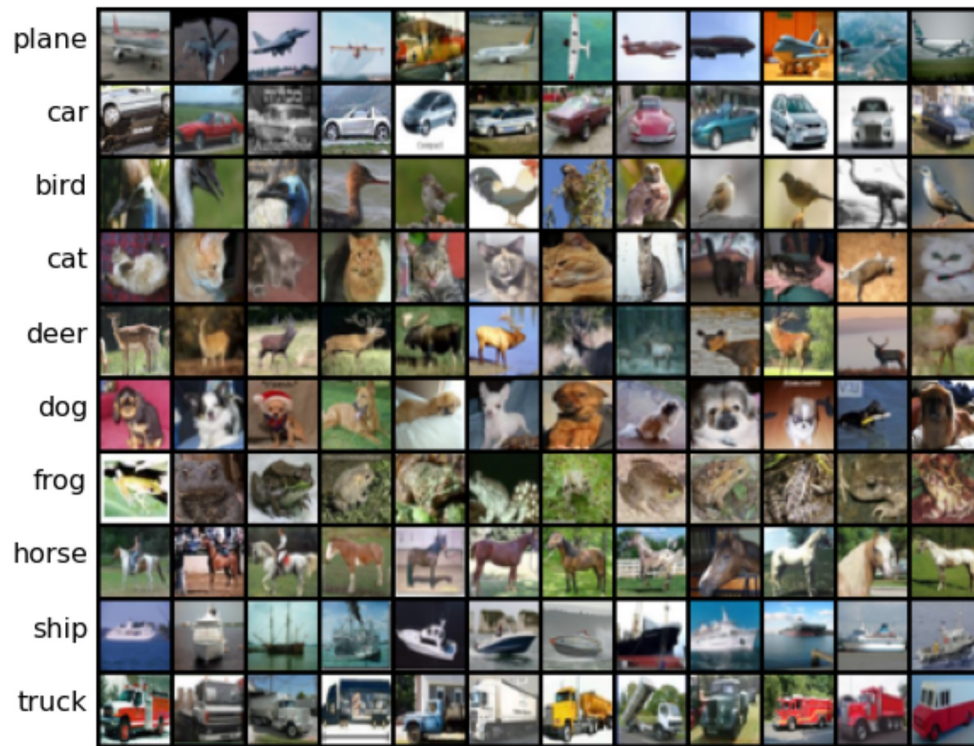
```
1 x_train, y_train, x_test, y_test = couils.data.cifar10()
2
3 print('Training set:', )
4 print(' data shape:', x_train.shape)
5 print(' labels shape: ', y_train.shape)
6 print('Test set:')
7 print(' data shape: ', x_test.shape)
8 print(' labels shape', y_test.shape)
```

```
Training set:
data shape: torch.Size([50000, 3, 32, 32])
labels shape: torch.Size([50000])
Test set:
data shape: torch.Size([10000, 3, 32, 32])
labels shape torch.Size([10000])
```

▼ Visualize the dataset

To give you a sense of the nature of the images in CIFAR-10, this cell visualizes some random examples from the training set.

```
1 import random
2 from torchvision.utils import make_grid
3
4 classes = ['plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
5 samples_per_class = 12
6 samples = []
7 for y, cls in enumerate(classes):
8     plt.text(-4, 34 * y + 18, cls, ha='right')
9     idxs = (y_train == y).nonzero().view(-1)
10    for i in range(samples_per_class):
11        idx = idxs[random.randrange(idxs.shape[0])].item()
12        samples.append(x_train[idx])
13 img = torchvision.utils.make_grid(samples, nrow=samples_per_class)
14 plt.imshow(couils.tensor_to_image(img))
15 plt.axis('off')
16 plt.show()
```



▼ Subsample the dataset

When implementing machine learning algorithms, it's usually a good idea to use a small sample of the full dataset. This way your code will run much faster, allowing for more interactive and efficient development. Once you are satisfied that you have correctly implemented the algorithm, you can then rerun with the entire dataset.

The function `coutils.data.cifar10()` can automatically subsample the CIFAR10 dataset for us. To see how to use it, we can check the documentation using the built-in `help` command:

```
1 help(coutils.data.cifar10)
```

```
Help on function cifar10 in module coutils.data:
```

```
cifar10(num_train=None, num_test=None)
    Return the CIFAR10 dataset, automatically downloading it if necessary.
    This function can also subsample the dataset.
```

Inputs:

- num_train: [Optional] How many samples to keep from the training set.
If not provided, then keep the entire training set.
- num_test: [Optional] How many samples to keep from the test set.
If not provided, then keep the entire test set.

Returns:

- x_train: float32 tensor of shape (num_train, 3, 32, 32)
- y_train: int64 tensor of shape (num_train, 3, 32, 32)
- x_test: float32 tensor of shape (num_test, 3, 32, 32)
- y_test: int64 tensor of shape (num_test, 3, 32, 32)

We will subsample the data to use only 500 training examples and 100 test examples:

```
1 num_train = 500
2 num_test = 250
3
4 x_train, y_train, x_test, y_test = cutils.data.cifar10(num_train, num_test)
5
6 print('Training set:', )
7 print(' data shape:', x_train.shape)
8 print(' labels shape: ', y_train.shape)
9 print('Test set:')
10 print(' data shape: ', x_test.shape)
11 print(' labels shape', y_test.shape)
```

Training set:

data shape: torch.Size([500, 3, 32, 32])
labels shape: torch.Size([500])

Test set:

data shape: torch.Size([250, 3, 32, 32])
labels shape torch.Size([250])

▼ Compute distances: Naive implementation

Now that we have examined and prepared our data, it is time to implement the kNN classifier. We can break the process down into two steps:

1. Compute the (squared Euclidean) distances between all training examples and all test examples
2. Given these distances, for each test example find its k nearest neighbors and have them vote for the label to output

Let's begin with computing the distance matrix between all training and test examples. First we will implement a naive version of the distance

Lets begin with computing the distance matrix between all training and test examples. First we will implement a naive version of the distance computation, using explicit loops over the training and test sets:

NOTE: When implementing distance functions in this notebook, you may not use the `torch.norm` function (or its instance method variant `x.norm`); you may not use any functions from `torch.nn` or `torch.nn.functional`.

```
1 def compute_distances_two_loops(x_train, x_test):
2     """
3     Computes the squared Euclidean distance between each element of the training
4     set and each element of the test set. Images should be flattened and treated
5     as vectors.
6
7     This implementation uses a naive set of nested loops over the training and
8     test data.
9
10    Inputs:
11    - x_train: Torch tensor of shape (num_train, C, H, W)
12    - x_test: Torch tensor of shape (num_test, C, H, W)
13
14    Returns:
15    - dists: Torch tensor of shape (num_train, num_test) where dists[i, j] is the
16        squared Euclidean distance between the ith training point and the jth test
17        point.
18    """
19    # Initialize dists to be a tensor of shape (num_train, num_test) with the
20    # same datatype and device as x_train
21    num_train = x_train.shape[0]
22    num_test = x_test.shape[0]
23    dists = x_train.new_zeros(num_train, num_test)
24    #####
25    # TODO: Implement this function using a pair of nested loops over the #
26    # training data and the test data. #
27    # #
28    # You may not use torch.norm (or its instance method variant), nor any #
29    # functions from torch.nn or torch.nn.functional. #
30    #####
31    for i in range(num_train):
32        for j in range(num_test):
33            dists[i, j] = (x_train[i, :, :, :] - x_test[j, :, :, :]).pow(2).sum().sqrt().item()
34    #####
35    #                               END OF YOUR CODE                               #
36    #####
37    return dists
```

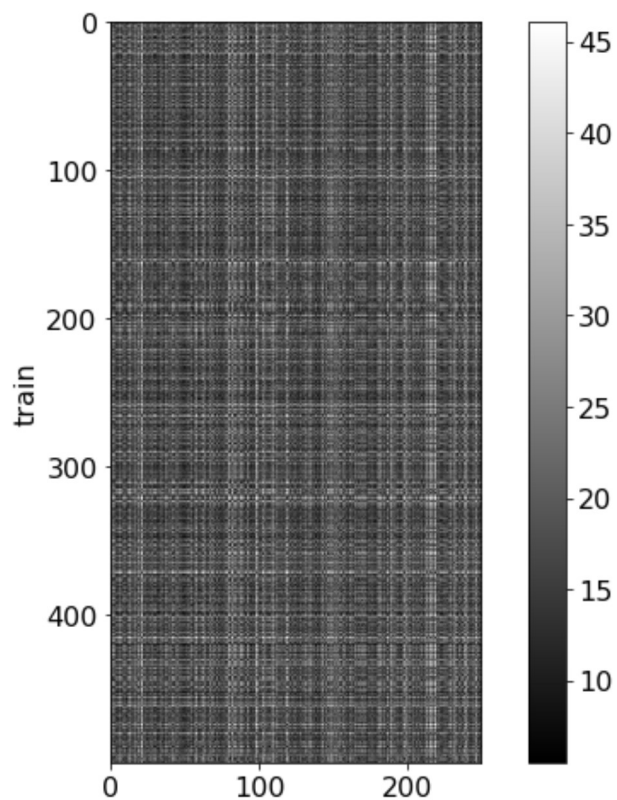
After implementing the function above, we can run it to check that it has the expected shape:

```
1 num_train = 500
2 num_test = 250
3 x_train, y_train, x_test, y_test = couils.data.cifar10(num_train, num_test)
4
5 dists = compute_distances_two_loops(x_train, x_test)
6 print('dists has shape: ', dists.shape)
```

```
dists has shape: torch.Size([500, 250])
```

As a visual debugging step, we can visualize the distance matrix, where each row is a test example and each column is a training example.

```
1 plt.imshow(dists.numpy(), cmap='gray', interpolation='none')
2 plt.colorbar()
3 plt.xlabel('test')
4 plt.ylabel('train')
5 plt.show()
```



test

▼ Compute distances: Vectorization

Our implementation of the distance computation above is fairly inefficient since it uses nested Python loops over the training and test sets.

When implementing algorithms in PyTorch, it's best to avoid loops in Python if possible. Instead it is preferable to implement your computation so that all loops happen inside PyTorch functions. This will usually be much faster than writing your own loops in Python, since PyTorch functions can be internally optimized to iterate efficiently, possibly using multiple threads. This is especially important when using a GPU to accelerate your code.

The process of eliminating explicit loops from your code is called **vectorization**. Sometimes it is straightforward to vectorize code originally written with loops; other times vectorizing requires thinking about the problem in a new way. We will use vectorization to improve the speed of our distance computation function.

As a first step toward vectorizing our distance computation, complete the following implementation which uses only a single Python loop over the training data:

```
1 def compute_distances_one_loop(x_train, x_test):
2     """
3     Computes the squared Euclidean distance between each element of the training
4     set and each element of the test set. Images should be flattened and treated
5     as vectors.
6
7     This implementation uses only a single loop over the training data.
8
9     Inputs:
10    - x_train: Torch tensor of shape (num_train, C, H, W)
11    - x_test: Torch tensor of shape (num_test, C, H, W)
12
13    Returns:
14    - dists: Torch tensor of shape (num_train, num_test) where dists[i, j] is the
15      squared Euclidean distance between the ith training point and the jth test
16      point.
17    """
18    # Initialize dists to be a tensor of shape (num_train, num_test) with the
19    # same datatype and device as x_train
20    num_train = x_train.shape[0]
21    num_test = x_test.shape[0]
22    dists = x_train.new_zeros(num_train, num_test)
23    #####
24    # TODO: Implement this function using only a single loop over x train.      #
```



```

25 #                                     #
26 # You may not use torch.norm (or its instance method variant), nor any      #
27 # functions from torch.nn or torch.nn.functional.                          #
28 #####
29 for i in range(num_train):
30     dists[i] = (x_test[:, :, :, :] - x_train[i, :, :, :]).pow(2).sum(dim=(1,2,3)).sqrt()
31 #####
32 #                                     #
33 #                                     #
34 return dists

```

We can check the correctness of our one-loop implementation by comparing it with our two-loop implementation on some randomly generated data.

Note that we do the comparison with 64-bit floating points for increased numeric precision.

```

1 torch.manual_seed(0)
2 x_train_rand = torch.randn(100, 3, 16, 16, dtype=torch.float64)
3 x_test_rand = torch.randn(100, 3, 16, 16, dtype=torch.float64)
4
5 dists_one = compute_distances_one_loop(x_train_rand, x_test_rand)
6 dists_two = compute_distances_two_loops(x_train_rand, x_test_rand)
7 difference = (dists_one - dists_two).pow(2).sum().sqrt().item()
8 print('Difference: ', difference)
9 if difference < 1e-4:
10     print('Good! The distance matrices match')
11 else:
12     print('Uh-oh! The distance matrices are different')

```

```

Difference:  0.0
Good! The distance matrices match

```

Now implement a fully vectorized version of the distance computation function that does not use any python loops.

```

1 def compute_distances_no_loops(x_train, x_test):
2     """
3     Computes the squared Euclidean distance between each element of the training
4     set and each element of the test set. Images should be flattened and treated
5     as vectors.
6
7     This implementation should not use any Python loops. For memory-efficiency,
8     it also should not create any large intermediate tensors; in particular you
9     should not create any intermediate tensors with O(num_train*num_test)

```

```

10 elements.
11
12 Inputs:
13 - x_train: Torch tensor of shape (num_train, C, H, W)
14 - x_test: Torch tensor of shape (num_test, C, H, W)
15
16 Returns:
17 - dists: Torch tensor of shape (num_train, num_test) where dists[i, j] is the
18   squared Euclidean distance between the ith training point and the jth test
19   point.
20 """
21 # Initialize dists to be a tensor of shape (num_train, num_test) with the
22 # same datatype and device as x_train
23 num_train = x_train.shape[0]
24 num_test = x_test.shape[0]
25 dists = x_train.new_zeros(num_train, num_test)
26 #####
27 # TODO: Implement this function without using any explicit loops and without #
28 # creating any intermediate tensors with O(num_train * num_test) elements. #
29 #                                                                           #
30 # You may not use torch.norm (or its instance method variant), nor any    #
31 # functions from torch.nn or torch.nn.functional.                         #
32 #                                                                           #
33 # HINT: Try to formulate the Euclidean distance using two broadcast sums    #
34 #   and a matrix multiply.                                                  #
35 #####
36 # Replace "pass" statement with your code
37 x_train = x_train.view(num_train, -1)
38 x_test = x_test.view(num_test, -1)
39 dists = (
40     -2 * x_train.mm(x_test.T) +
41     x_test.pow(2).sum(dim=1).view(1, -1) +
42     x_train.pow(2).sum(dim=1).view(-1, 1)
43 ).sqrt()
44
45 """
46 # Another version:
47 dists = (
48     (x_test.view(1, num_test, -1) - x_train.view(num_train, 1, -1))
49     .pow(2)
50     .sum(dim=(2))
51     .sqrt()
52 )
53 # But it fails with:
54 # RuntimeError: CUDA out of memory. Tried to allocate 28.61 GiB (GPU 0; 15.90 GiB total capacity;
55 # 74.03 MiB already allocated; 15.20 GiB free; 82.00 MiB reserved in total by PyTorch)
56 """
57 #####
58 #                               END OF YOUR CODE                               #

```

```

58 # ##### END OF YOUR CODE #####
59 #####
60 return dists

```

As before, we can check the correctness of our implementation by comparing the fully vectorized version against the original naive version:

```

1 torch.manual_seed(0)
2 x_train_rand = torch.randn(100, 3, 16, 16, dtype=torch.float64)
3 x_test_rand = torch.randn(100, 3, 16, 16, dtype=torch.float64)
4
5 dists_two = compute_distances_two_loops(x_train_rand, x_test_rand)
6 dists_none = compute_distances_no_loops(x_train_rand, x_test_rand)
7 difference = (dists_two - dists_none).pow(2).sum().sqrt().item()
8 print('Difference: ', difference)
9 if difference < 1e-4:
10     print('Good! The distance matrices match')
11 else:
12     print('Uh-oh! The distance matrices are different')

```

```

Difference:  3.1760544968942397e-13
Good! The distance matrices match

```

We can now compare the speed of our three implementations. If you've implemented everything properly, the one-loop implementation should take less than 4 seconds to run, and the fully vectorized implementation should take less than 0.1 seconds to run.

```

1 import time
2
3 def timeit(f, *args):
4     tic = time.time()
5     f(*args)
6     toc = time.time()
7     return toc - tic
8
9 torch.manual_seed(0)
10 x_train_rand = torch.randn(500, 3, 32, 32)
11 x_test_rand = torch.randn(500, 3, 32, 32)
12
13 two_loop_time = timeit(compute_distances_two_loops, x_train_rand, x_test_rand)
14 print('Two loop version took %.2f seconds' % two_loop_time)
15
16 one_loop_time = timeit(compute_distances_one_loop, x_train_rand, x_test_rand)
17 speedup = two_loop_time / one_loop_time
18 print('One loop version took %.2f seconds (%.1fX speedup)'
19       % (one_loop_time, speedup))

```

```

20
21 no_loop_time = timeit(compute_distances_no_loops, x_train_rand, x_test_rand)
22 speedup = two_loop_time / no_loop_time
23 print('No loop version took %.2f seconds (%.1fX speedup)'
24       % (no_loop_time, speedup))

```

```

Two loop version took 9.26 seconds
One loop version took 0.66 seconds (13.9X speedup)
No loop version took 0.02 seconds (555.5X speedup)

```

▼ Predict labels

Now that we have a method for computing distances between training and test examples, we need to implement a function that uses those distances together with the training labels to predict labels for test samples.

Complete the implementation of the function below:

```

1 def predict_labels(dists, y_train, k=1):
2     """
3     Given distances between all pairs of training and test samples, predict a
4     label for each test sample by taking a majority vote among its k nearest
5     neighbors in the training set.
6
7     In the event of a tie, this function should return the smaller label. For
8     example, if k=5 and the 5 nearest neighbors to a test example have labels
9     [1, 2, 1, 2, 3] then there is a tie between 1 and 2 (each have 2 votes), so
10    we should return 1 since it is the smaller label.
11
12    Inputs:
13    - dists: Torch tensor of shape (num_train, num_test) where dists[i, j] is the
14      squared Euclidean distance between the ith training point and the jth test
15      point.
16    - y_train: Torch tensor of shape (y_train,) giving labels for all training
17      samples. Each label is an integer in the range [0, num_classes - 1]
18    - k: The number of nearest neighbors to use for classification.
19
20    Returns:
21    - y_pred: A torch int64 tensor of shape (num_test,) giving predicted labels
22      for the test data, where y_pred[j] is the predicted label for the jth test
23      example. Each label should be an integer in the range [0, num_classes - 1].
24    """
25    num_train, num_test = dists.shape
26    y_pred = torch.zeros(num_test, dtype=torch.int64)
27    #####
28    # TODO: Implement this function. You may use an explicit loop over the test #

```

```

28 # TODO: Implement this function. You may use an explicit loop over the test #
29 # samples. Hint: Look up the function torch.topk #
30 #####
31 # Replace "pass" statement with your code
32 for j in range(num_test):
33     y_pred[j] = y_train[
34         torch.topk(dists[:,j], k=k, largest=False, sorted=True).indices
35         ].bincount().argmax(dim=0)
36
37 #####
38 #                               END OF YOUR CODE                               #
39 #####
40 return y_pred

```

Now we have implemented all the required functionality for the K-Nearest Neighbor classifier. We can define a simple object to encapsulate the classifier:

```

1 class KnnClassifier:
2     def __init__(self, x_train, y_train):
3         """
4         Create a new K-Nearest Neighbor classifier with the specified training data.
5         In the initializer we simply memorize the provided training data.
6
7         Inputs:
8         - x_train: Torch tensor of shape (num_train, C, H, W) giving training data
9         - y_train: int64 torch tensor of shape (num_train,) giving training labels
10        """
11        self.x_train = x_train.contiguous()
12        self.y_train = y_train.contiguous()
13
14    def predict(self, x_test, k=1):
15        """
16        Make predictions using the classifier.
17
18        Inputs:
19        - x_test: Torch tensor of shape (num_test, C, H, W) giving test samples
20        - k: The number of neighbors to use for predictions
21
22        Returns:
23        - y_test_pred: Torch tensor of shape (num_test,) giving predicted labels
24          for the test samples.
25        """
26        dists = compute_distances_no_loops(self.x_train, x_test.contiguous())
27        y_test_pred = predict_labels(dists, self.y_train, k=k)
28        return y_test_pred
29

```

```

30 def check_accuracy(self, x_test, y_test, k=1, quiet=False):
31     """
32     Utility method for checking the accuracy of this classifier on test data.
33     Returns the accuracy of the classifier on the test data, and also prints a
34     message giving the accuracy.
35
36     Inputs:
37     - x_test: Torch tensor of shape (num_test, C, H, W) giving test samples
38     - y_test: int64 torch tensor of shape (num_test,) giving test labels
39     - k: The number of neighbors to use for prediction
40     - quiet: If True, don't print a message.
41
42     Returns:
43     - accuracy: Accuracy of this classifier on the test data, as a percent.
44       Python float in the range [0, 100]
45     """
46     y_test_pred = self.predict(x_test, k=k)
47     num_samples = x_test.shape[0]
48     num_correct = (y_test == y_test_pred).sum().item()
49     accuracy = 100.0 * num_correct / num_samples
50     msg = (f'Got {num_correct} / {num_samples} correct; '
51           f'accuracy is {accuracy:.2f}%')
52     if not quiet:
53         print(msg)
54     return accuracy

```

Now lets put everything together and test our K-NN classifier on a subset of CIFAR-10, using k=1:

If you've implemented everything correctly you should see an accuracy of about 27%.

```

1 num_train = 5000
2 num_test = 500
3 x_train, y_train, x_test, y_test = cutils.data.cifar10(num_train, num_test)
4 classifier = KnnClassifier(x_train, y_train)
5 classifier.check_accuracy(x_test, y_test, k=1)

```

```

Got 137 / 500 correct; accuracy is 27.4%
27.4

```

Now lets increase to k=5. You should see a slightly higher accuracy than k=1:

```

1 classifier.check_accuracy(x_test, y_test, k=5)

```

```

Got 139 / 500 correct; accuracy is 27.80%

```

▼ Cross-validation

We have not implemented the full k-Nearest Neighbor classifier, but the choice of $k = 5$ was arbitrary. We will use **cross-validation** to set this hyperparameter in a more principled manner.

Implement the following function to run cross-validation:

```

1 def knn_cross_validate(x_train, y_train, num_folds=5, k_choices=None):
2     """
3     Perform cross-validation for KnnClassifier.
4
5     Inputs:
6     - x_train: Tensor of shape (num_train, C, H, W) giving all training data
7     - y_train: int64 tensor of shape (num_train,) giving labels for training data
8     - num_folds: Integer giving the number of folds to use
9     - k_choices: List of integers giving the values of k to try
10
11    Returns:
12    - k_to_accuracies: Dictionary mapping values of k to lists, where
13      k_to_accuracies[k][i] is the accuracy on the ith fold of a KnnClassifier
14      that uses k nearest neighbors.
15    """
16    if k_choices is None:
17        # Use default values
18        k_choices = [1, 3, 5, 8, 10, 12, 15, 20, 50, 100]
19
20    # First we divide the training data into num_folds equally-sized folds.
21    x_train_folds = []
22    y_train_folds = []
23    #####
24    # TODO: Split the training data and images into folds. After splitting,      #
25    # x_train_folds and y_train_folds should be lists of length num_folds, where #
26    # y_train_folds[i] is the label vector for images in x_train_folds[i].      #
27    # Hint: torch.chunk                                                         #
28    #####
29    # Replace "pass" statement with your code
30    x_train_folds = x_train.chunk(num_folds)
31    y_train_folds = y_train.chunk(num_folds)
32    #####
33    #                                     END OF YOUR CODE                       #
34    #####
35
36    # A dictionary holding the accuracies for different values of k that we find

```

```

37 # when running cross-validation. After running cross-validation,
38 # k_to_accuracies[k] should be a list of length num_folds giving the different
39 # accuracies we found when trying KnnClassifiers that use k neighbors.
40 k_to_accuracies = {}
41
42 #####
43 # TODO: Perform cross-validation to find the best value of k. For each value #
44 # of k in k_choices, run the k-nearest-neighbor algorithm num_folds times; #
45 # in each case you'll use all but one fold as training data, and use the #
46 # last fold as a validation set. Store the accuracies for all folds and all #
47 # values in k in k_to_accuracies. HINT: torch.cat #
48 #####
49 # Replace "pass" statement with your code
50 for k in k_choices:
51     for fi in range(num_folds):
52         f_x_train = torch.cat(x_train_folds[:fi] + x_train_folds[fi+1:], dim=0)
53         f_y_train = torch.cat(y_train_folds[:fi] + y_train_folds[fi+1:], dim=0)
54         classifier = KnnClassifier(f_x_train, f_y_train)
55         acc = classifier.check_accuracy(x_train_folds[fi], y_train_folds[fi], k=k, quiet=True)
56         if k not in k_to_accuracies:
57             k_to_accuracies[k] = []
58             k_to_accuracies[k] += [acc]
59 #####
60 #                               END OF YOUR CODE                               #
61 #####
62
63 return k_to_accuracies

```

Now we'll run your cross-validation function:

```

1 num_train = 5000
2 num_test = 500
3 x_train, y_train, x_test, y_test = couutils.data.cifar10(num_train, num_test)
4
5 k_to_accuracies = knn_cross_validate(x_train, y_train, num_folds=5)
6
7 for k, accs in sorted(k_to_accuracies.items()):
8     print('k = %d got accuracies: %r' % (k, accs))

```

```

k = 1 got accuracies: [26.3, 25.7, 26.4, 27.8, 26.6]
k = 3 got accuracies: [23.9, 24.9, 24.0, 26.6, 25.4]
k = 5 got accuracies: [24.8, 26.6, 28.0, 29.2, 28.0]
k = 8 got accuracies: [26.2, 28.2, 27.3, 29.0, 27.3]
k = 10 got accuracies: [26.5, 29.6, 27.6, 28.4, 28.0]
k = 12 got accuracies: [26.0, 29.5, 27.9, 28.3, 28.0]
k = 15 got accuracies: [25.2, 28.9, 27.8, 28.2, 27.4]

```



```

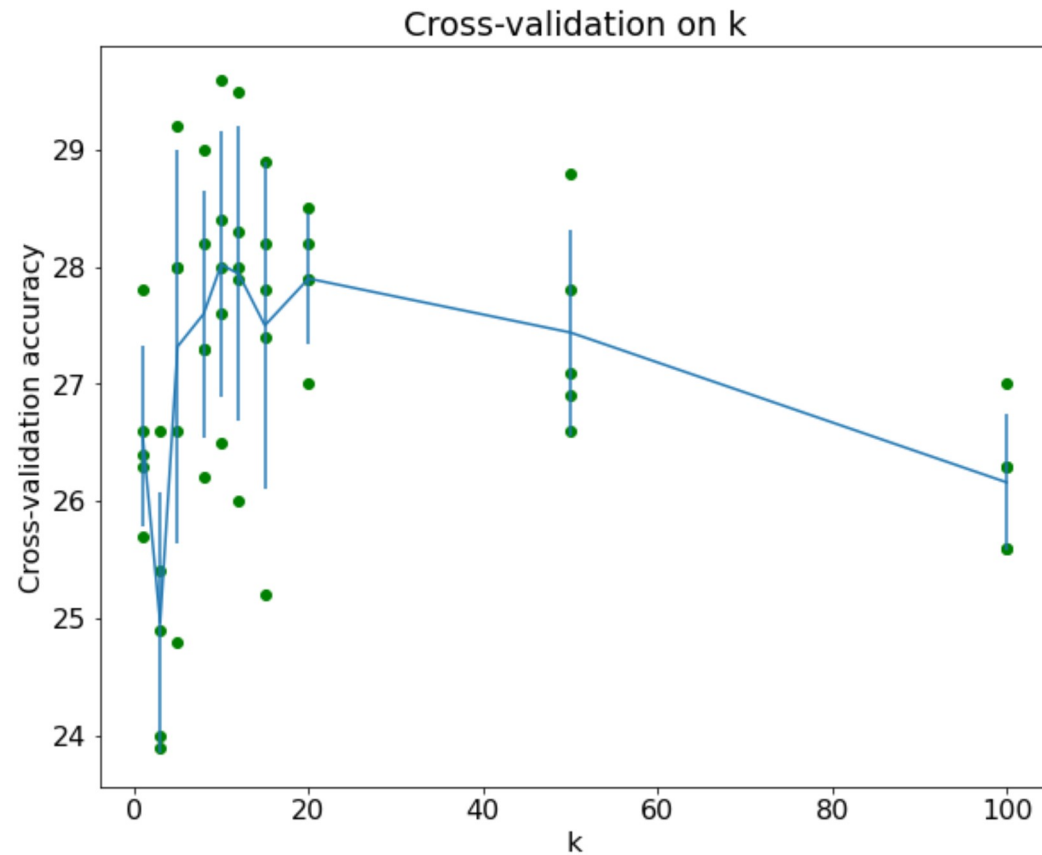
k = 20 got accuracies: [27.0, 27.9, 27.9, 28.2, 28.5]
k = 50 got accuracies: [27.1, 28.8, 27.8, 26.9, 26.6]
k = 100 got accuracies: [25.6, 27.0, 26.3, 25.6, 26.3]

```

```

1 ks, means, stds = [], [], []
2 for k, accs in sorted(k_to_accuracies.items()):
3     plt.scatter([k] * len(accs), accs, color='g')
4     ks.append(k)
5     means.append(statistics.mean(accs))
6     stds.append(statistics.stdev(accs))
7 plt.errorbar(ks, means, yerr=stds)
8 plt.xlabel('k')
9 plt.ylabel('Cross-validation accuracy')
10 plt.title('Cross-validation on k')
11 plt.show()

```



Now we can use the results of cross-validation to select the best value for k, and rerun the classifier on our full 5000 set of training examples.

You should get an accuracy above 28%.

```
1 best_k = 1
2 #####
3 # TODO: Use the results of cross-validation stored in k_to_accuracies to      #
4 # choose the value of k, and store the result in best_k. You should choose  #
5 # the value of k that has the highest mean accuracy accross all folds.      #
6 #####
7 # Replace "pass" statement with your code
8 best_k = max((statistics.mean(k_to_accuracies[k]), k) for k in k_to_accuracies.keys())[1]
9 #####
10 #                                END OF YOUR CODE                                #
11 #####
12
13 print('Best k is ', best_k)
14 classifier = KnnClassifier(x_train, y_train)
15 classifier.check_accuracy(x_test, y_test, k=best_k)
```

```
Best k is 10
Got 141 / 500 correct; accuracy is 28.20%
28.2
```

Finally, we can use our chosen value of k to run on the entire training and test sets.

This may take a while to run, since the full training and test sets have 50k and 10k examples respectively. You should get an accuracy above 33%.

Run this only once!

```
1 x_train_all, y_train_all, x_test_all, y_test_all = coutils.data.cifar10()
2 classifier = KnnClassifier(x_train_all, y_train_all)
3 classifier.check_accuracy(x_test_all, y_test_all, k=best_k)
```

```
Got 3386 / 10000 correct; accuracy is 33.86%
33.86
```

