

knn

January 20, 2024

1 Assignment 1: KNN

For this part of assignment, you are tasked to implement KNN algorithm and test it on the a subset of CIFAR10 dataset.

You could run the whole notebook and answer the question in the notebook.

TO SUBMIT: PDF of this notebook with all the required outputs and answers.

```
[11]: # Import Packages
import numpy as np
import matplotlib.pyplot as plt
```

1.1 Prepare Dataset

Since CIFAR10 is a relative large dataset, and KNN is quite time-consuming method, we only a small sub-set of CIFAR10 for KNN part

```
[12]: from utils.data_processing import get_cifar10_data

# Use a subset of CIFAR10 for KNN assignments
dataset = get_cifar10_data(subset_train=5000, subset_val=250, subset_test=500)

print(dataset.keys())
print("Training Set Data Shape: ", dataset["x_train"].shape)
print("Training Set Label Shape: ", dataset["y_train"].shape)
```

```
dict_keys(['x_train', 'y_train', 'x_val', 'y_val', 'x_test', 'y_test'])
Training Set Data Shape: (5000, 3072)
Training Set Label Shape: (5000,)
```

1.2 Implementation (60%)

You need to implement the KNN method in `algorithms/knn.py`. You need to fill in the prediction function(since the training of KNN is just remembering the training set).

For KNN implementation, you are tasked to implement two version of it.

- Two Loop Version: use one loop to iterate through training samples and one loop to iterate through test samples

- One Loop Version: use one loop to iterate through test samples and use broadcast (<https://numpy.org/doc/stable/user/basics.broadcasting.html>) feature of numpy to calculate all the distance at once

Note: It is possible to build a Fully Vectorized Version without explicit for loop to calculate the distance, but you do not have to do it in this assignment. You could use the fully vectorized version to replace the loop versions as well.

For distance function, in this assignment, we use Euclidean distance between samples.

```
[13]: from algorithms import KNN

knn = KNN(num_class=10)
knn.train(
    x_train=dataset["x_train"],
    y_train=dataset["y_train"],
    k=5,
)
```

1.2.1 Compare the time consumption of different method

In this section, you will test your different implementation of KNN method, and compare their speed.

```
[14]: from utils.evaluation import get_classification_accuracy
```

Two Loop Version:

```
[15]: import time

c_t = time.time()
prediction = knn.predict(dataset["x_test"], loop_count=2)
print("Two Loop Prediction Time:", time.time() - c_t)

test_acc = get_classification_accuracy(prediction, dataset["y_test"])
print("Test Accuracy:", test_acc)
```

Two Loop Prediction Time: 12.961191892623901
Test Accuracy: 0.278

One Loop Version

```
[16]: import time

c_t = time.time()
prediction = knn.predict(dataset["x_test"], loop_count=1)
print("One Loop Prediction Time:", time.time() - c_t)

test_acc = get_classification_accuracy(prediction, dataset["y_test"])
print("Test Accuracy:", test_acc)
```

One Loop Prediction Time: 14.301531076431274

Test Accuracy: 0.278

Your different implementation should output the exact same result

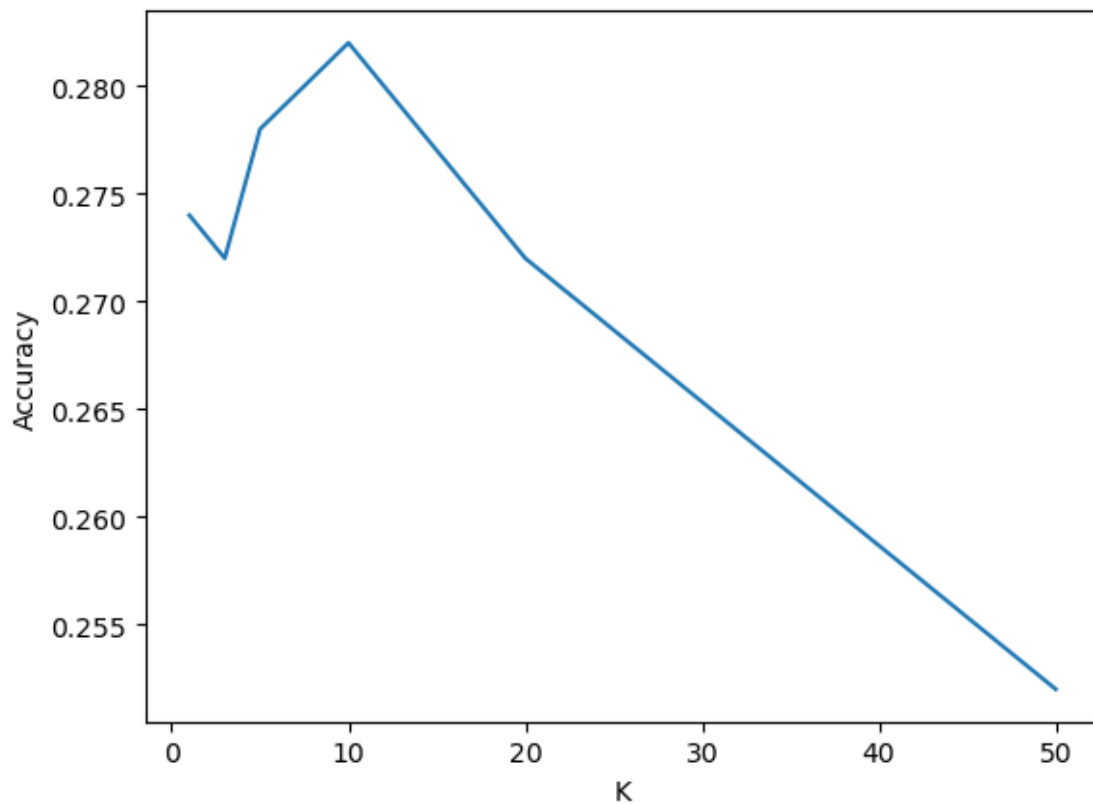
1.3 Test different Hyper-parameter (20%)

For KNN, there is only one hyper-parameter of the algorithm: How many nearest neighbour to use(**K**).

Here, you are provided the code to test different k for the same dataset.

```
[17]: accuracies = []

k_candidates = [1, 3, 5, 10, 20, 50]
for k_cand in k_candidates:
    prediction = knn.predict(x_test=dataset["x_test"], k=k_cand)
    acc = get_classification_accuracy(prediction, dataset["y_test"])
    accuracies.append(acc)
plt.ylabel("Accuracy")
plt.xlabel("K")
plt.plot(k_candidates, accuracies)
plt.show()
```



1.3.1 Inline Question 1:

Please describe the output result you get, and provide some explanation as well.

1.3.2 Your Answer:

It has the highest accuracy when $k = 12$. It starts decreasing when k is getting larger because k is larger than the actual number of categories in the dataset. The accuracy was increasing from $k = 0$ to $k = 12$ because k is smaller than the number of categories in the dataset and start increasing to the actual number of categories of the dataset.

1.4 Try different feature representation (19%)

Since machine learning method rely heavily on the feature extraction, you will see how different feature representation affect the performance of the algorithm in this section.

You are provided the code about using **HOG** descriptor to represent samples in the notebook.

```
[18]: from utils.data_processing import get_cifar10_data
      from utils.data_processing import HOG_preprocess
      from functools import partial

      # Delete previous dataset to save memory
      del dataset
      del knn

      # Use a subset of CIFAR10 for KNN assignments
      hog_p_func = partial(
          HOG_preprocess,
          orientations=9,
          pixels_per_cell=(4, 4),
          cells_per_block=(1, 1),
          visualize=False,
          # multichannel=True,
          channel_axis=2
      )
      dataset = get_cifar10_data(
          feature_process=hog_p_func, subset_train=5000, subset_val=250,
          ↪subset_test=500
      )
```

Start Processing

Processing Time: 2.981959819793701

```
[19]: knn = KNN(num_class=10)
      knn.train(
          x_train=dataset["x_train"],
          y_train=dataset["y_train"],
          k=5,
```

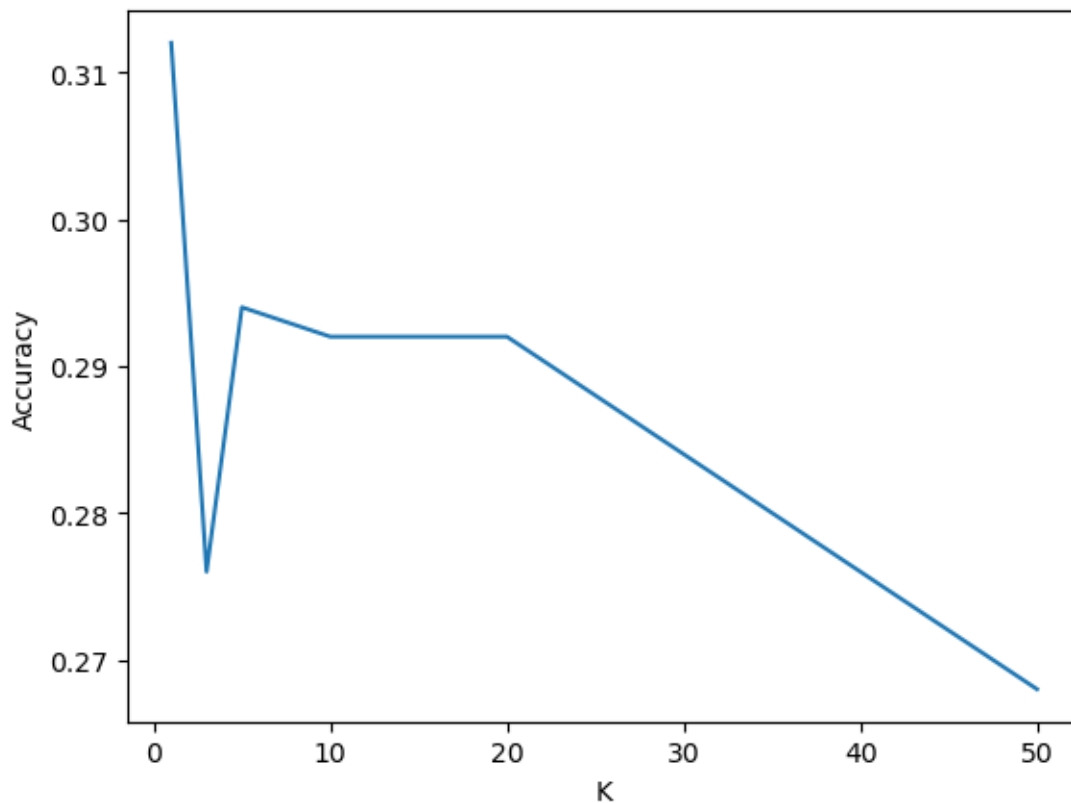
```

)
accuracies = []

k_candidates = [1, 3, 5, 10, 20, 50]
for k_cand in k_candidates:
    prediction = knn.predict(x_test=dataset["x_test"], k=k_cand)
    acc = get_classification_accuracy(prediction, dataset["y_test"])
    accuracies.append(acc)

plt.ylabel("Accuracy")
plt.xlabel("K")
plt.plot(k_candidates, accuracies)
plt.show()

```



1.4.1 Inline Question 2:

Please describe the output result you get, compare with the result you get in the previous section, and provide some explanation as well.

1.4.2 Your Answer:

The result I got here is different from the previous(Higher accuracy and different k value for highest accuracy) because we use different feature extraction method to process the dataset and machine learning performance is largely depend on feature extraction.

1.5 Survey (1%)

1.5.1 Question:

How many hours did you spend on assignment 1?

1.5.2 Your Answer: I spend about 6 hours on completing assignment 1.

[]: