

Proposal Report On Parallel Implementation of K-NN Classifier

Kajol Jain
NITK
Computer Science Engineering
Roll No- 17CS014
Email: kajoljain379@gmail.com

Priyanka Gupta
NITK
Computer Science Engineering
Roll No- 17CS019
Email: prigupta9875@gmail.com

I. PROBLEM STATEMENT

k-Nearest Neighbor (k-NN) is a classification algorithm used in machine learning and data mining applications such as email spam filtering, content retrieval, customer segmentation in online shopping websites etc. In the algorithm, the variable k is defined by the user. k-NN classifiers find the k number for training set, which has similar or is closest to the test data. The k-NN algorithm can be used for classification, based on the distance to the k nearest members, the algorithm decides which class the given input should belong to.

II. OBJECTIVES

High parallelism can be achieved using GPU and in comparatively lesser cost than CPU. So we are trying to do the same task in GPU which speeds up the execution by 99%. We will do the implementation consists of different level of parallelism like searching, sorting and other parallelly executable task.

III. WORK DISTRIBUTION AMONG TEAM MEMBERS

This algorithm has three main steps and they are,
1) distance calculation 2) nearest neighbor search 3) class prediction
So, we will apply the parallelism on each above steps and each teammate will work on one algorithm and together merge it with class prediction step.

IV. IMPLEMENTATION

We have done the serial implementation of KNN algorithm which focuses on how to classify a large dataset which takes significant amount of time in a conventional CPU. We have done this via following steps i.e.

A. Distance calculation:

We have considered the SHUTTLE numeric dataset from the website of Machine Learning Repository. The dataset consists of 43,500 training data and 14,500 testing data, each having nine attributes without any missing components. The distance calculation should measure how close the given vector (also known as the test vector) is to all of the training vectors. We used the Euclidian distance (D) method to calculate

the distance between two vectors X and Y of size n can be calculated using Equation 1.

$$D = \sqrt{\sum_{i=1}^n (X_i - Y_i)^2} \quad (1)$$

B. Nearest Neighbor Search:

In the second step, minimum k distance values will be sorted with their corresponding class values for the test data. So, these k values can be used to predict final class.

C. Class prediction:

The final phase locates the majority class value in the sorted training set and predicts that class to the test data. So, occurrence of each class is calculated that whichever class occur max, that class is taken as corresponding class for that point and this procedure for all test data.

V. RESULT OF SERIAL IMPLEMENTATION

We implement Serial algorithm of KNN classifier with all steps mentioned on dataset taken from website Machine learning repository. Here we run each step serially on 14500 test data by calculating its Euclidian distance from all 43500 train dataset. We also check our result from given data result and it predicts resultant class correctly for 99.9% of test data. We check it using different values of k.

VI. PROPOSED IMPROVEMENTS

In referred paper by us, they apply sorting on full distance array but as we required only k nearest neighbours so, in our implementation we apply selection sort k times which improve execution time from $O(n \log n)$ to $O(n)$.

VII. PARALLEL IMPLEMENTATION

KNN algorithm is data parallel application, so we can improve its performance drastically using parallel computing which improves it near by 100 times in compare to serial implementation on CPU. So, there we make steps parallel like distance calculation and Nearest neighbour classification. This serial implementation with k value 12 takes 43 second on CPU.

A. Distance calculation

We designed the process as each block responsible for calculating the distance between one test and all training data set. This kernel returns all the distances as a 1D array to find the k-nearest neighbors. The length of the array is the multiplication of the number of test records and the number of training records.

B. Nearest Neighbour

In this kernel, we define the number of threads as equal to the number of testing records in the testing dataset. We send only the distance array one test data per thread. It finds the index of minimum k distance for each test data.

C. Class Prediction

We find the corresponding class values from class array and replace the index with the corresponding class value and we return the k size array for each test record. This kernel returns k-dimensional array for each test record to class prediction function which is implemented in CPU.

VIII. EXPERIMENTAL SETUP

The CPU program is tested on a machine with Intel core i5-3470 CPU @ 3.20GHz processor and 8GB memory. The GPU program is tested on Tesla K40, which consists of 2880 cores and supports up to 30720 threads. The GPU memory was 12GB. We have considered the SHUTTLE numeric dataset from the website of Machine Learning Repository. The dataset consists of 43,500 training data and 14,500 testing data, each having nine attributes without any missing components. The dataset was stored in the CSV format. We have divided the dataset into various sizes of training sets and test sets to evaluate the execution time of the distance calculation phase and the nearest neighbor search of the k-NN algorithm across GPU and CPU. The number of neighbors (k value) is taken as 5 for all the test cases.

A. First Implementation

First implementation was the serial implementation as we discussed above. In this we calculate the Euclidian distance of 14500 test point with all the train data i.e. 43500 trained data points in a serial manner. Each data point has 9 attributes. After finding the distance we extract the nearest k value from the distance array and then find out the nearest neighbour class. This serial computation on CPU takes more time whenever data increases fastly. To minimize the time taken by serial implementation we implement the parallel version of implementation by using GPU supported architecture. When the test data size increases, the execution time of both CPU and GPU increases linearly. When the test data size is 14,500, the CPU takes about 43 seconds.

B. Second Implementation

We implement the parallel program of K-NN classifier. We write the two kernel modules: one for calculating the Euclidian distance and other for finding k nearest neighbour. While calculating the distance we designed the process as each block responsible for calculating the distance between one test and all training data set. Inside the block, each thread calculates the distance between one test data and one training data. The blocks are arranged in a 2D thread organization. The corresponding attributes between testing and training dataset are subtracted first. After that, it takes the square root of summation of subtracted values and stores that value in a correct index of the distance array. Index identifying like $index = TrainRecordIndex * NumberOfTestRecord + TestRecordIndex$.

This kernel returns all the distances as a 1D array to find the k-nearest neighbors. The length of the array is the multiplication of the number of test records and the number of training records.

This array will pass to the other kernel i.e. K-NN nearest neighbour for predicting the corresponding class. First it sorts the first k nearest distance and then finds out that from which class majority of data points belong that will be the resultant class of the corresponding test point. This will take the approx 20 sec for more than 14500 test point with 9 attributes. This implementation will speed up the execution more than 70 %.

C. Third Implementation

We have presented a GPU based implementation of k-NN classification algorithm with two kernels, the distances calculation kernel and the nearest neighbor search kernel. There are some limitations in the current implementation of k-NN algorithm in GPU. Though the program results for a large and common machine learning dataset, but the program fails to handle larger data, i.e. when it exceeds the device hardware memory and loaded into device memory. The reason is that the program tries to do all the work at once loading the whole training and testing data set into the memory.

Therefore, we try to develop the chunk concept. We divide the test data in the form of chunk of small size according to device hardware memory available.

We are doing the Euclidian distance calculation and k nearest neighbour parallelly but the chunkwise. After that we store the result of each chunk which is the storing of k nearest neighbour of each chunk data and at last we combine all the data of each chunk and at last we do class prediction of all the stored data combinely. We use 100000 training data and 40000 test data in chunk implementation. It will take 135 sec in serial implementation and 56 sec in chunk implementation. Which is not possible with these much data in our first parallel implementation due to device hardware implementation because of no of thread created by these much data.

```

1. __global__ void euclideanDistance() {
2.
3.   int row = blockIdx.x * blockDim.x +
   threadIdx.x;
4.   int col = blockIdx.y * blockDim.y + threadIdx.y
5.   distanceId = row * testRecords+col;
6.
7.   if (row less than trainRecords and col less
   than
8.   testRecords)
9.   {
10.    for (i = 0 to attributes - 1) {
11.     diff <= (d_trainSet [row*attributes + i] -
12.     d_testSet[col*attributes +
13.     i])
14.     sum += diff * diff;
15.    }
16.    distanceArray [distanceId] <= sqrt(sum)
17.  }
18. }

```

Fig .1: Kernel code for Euclidian distance

```

1. __global__ void NearestNeighbour(){
2.
3.   int threadId = blockIdx.x * blockDim.x +
   threadIdx.x;
4.
5.   if(threadId less than testRecords){
6.     for(b = 0 to < k){
7.       Min <= distanceArray[threadId]
8.       index = 0;
9.       for(c = 1 to c){
10.        if((Min greater than distanceArray[c*testRecords
11.        + threadId]){
12.          Min <= distanceArray[c*testRecords + threadId]
13.          index <= c
14.        }
15.        if((Max less than distanceArray[c*testRecords +
16.        threadId] and (b equal to 0)){
17.          Max <= distanceArray[c*testRecords + threadId]
18.        }
19.        distanceArray[index*testRecords + threadId] <=
20.        Max
21.        class[testRecords*b +threadId] <=
22.        trainClass[index]
23.      }
24.    }
25.  }
26. }

```

Fig.2 : kernel code for nearest neighbour

```

1. __global__
2. void class_classification(int *index_chunks,double **set, int
   total_chunks_train ,int k,int *d_kneighbours,int set1,int
   *res_class){
3.   int i=blockDim.x*blockIdx.x+threadIdx.x;
4.   int set_i;
5.   int min=0;
6.   for(int f=0;f<total_chunks_train;f++)
7.   index_chunks[f]=0;
8.   if(i<test_row)
9.   {
10.    for(int l=0;l<k;l++)
11.    {
12.     min=0;
13.     for(int j=1;j<total_chunks_train;j++)
14.     {
15.      if(set[min][((index_chunks[min]*test_row+i)*2)]>set[j][
16.      ((index_chunks[j]*test_row+i)*2)])
17.      min=j;
18.     }
19.     set_i=i*set1;
20.     d_kneighbours[(int)(set[min][((index_chunks[min]*test_row+i)*2
21.     +1))-1+set_i]+=1;
22.     index_chunks[min]++;
23.    }
24.    set_i=i*set1;
25.    int max=0;
26.    for(int l=1;l<set1;l++)
27.    {
28.     if(d_kneighbours[set_i+l]>d_kneighbours[set_i+max])
29.     max=l;
30.    }
31.    res_class[i]=max+1;
32.  }
33. }

```

Fig.3 : kernel code for nearest neighbour in chunk implementation

X. CONCLUSION

We have presented three implementation first is serial implementation ,second is a GPU based implementation of k-NN classification algorithm with two kernels, the distances calculation kernel and the nearest neighbor search kernel.

We tested the solution using a standard dataset in both CPU and GPU and compared the results separately. We could find a significant speedup of the algorithm specially when the test data size becomes larger. For a test data set of 14,000, we could observe around 70x speedup in GPU.In the third implementation we tried to remove the device memory hardware limitation by introducing the chunk concept.So that we can the big data with the same hardware