## a5/cuda_kmeans_naive.cu

```c
1   #include <stdio.h>
2   #include <stdlib.h>
3
4   #include "kmeans.h"
5   #include "alloc.h"
6   #include "error.h"
7
8   #ifdef __CUDACC__
9   inline void checkCuda(cudaError_t e) {
10      if (e != cudaSuccess) {
11          // cudaGetErrorString() isn't always very helpful. Look up the error
12          // number in the cudaError enum in driver_types.h in the CUDA includes
13          // directory for a better explanation.
14          error("CUDA Error %d: %s\n", e, cudaGetErrorString(e));
15      }
16  }
17
18  inline void checkLastCudaError() {
19      checkCuda(cudaGetLastError());
20  }
21  #endif
22
23  __device__ int get_tid() {
24    return blockIdx.x * blockDim.x + threadIdx.x;
25  }
26
27  /* square of Euclid distance between two multi-dimensional points */
28  __host__ __device__ inline static
29  double euclid_dist_2(int numCoords,
30                       int numObjs,
31                       int numClusters,
32                       double *objects,     // [numObjs][numCoords]
33                       double *clusters,    // [numClusters][numCoords]
34                       int objectId,
35                       int clusterId) {
36    int i;
37    double ans = 0.0;
38
39    /* TODO: Calculate the euclid_dist of elem=objectId of objects from elem=clusterId from
    clusters*/
40    for (i = 0; i < numCoords; i++) {
41      double objectVal = objects[objectId * numCoords + i];
42      double clusterVal = clusters[clusterId * numCoords + i];
43
44      double diff = objectVal - clusterVal;
45      ans += diff * diff;
46    }
47
48    return (ans);
49  }
50
51  __global__ static
```

```
52  void find_nearest_cluster(int numCoords,
53                            int numObjs,
54                            int numClusters,
55                            double *objects,         //  [numObjs][numCoords]
56                            double *deviceClusters,  //  [numClusters][numCoords]
57                            int *deviceMembership,      //  [numObjs]
58                            double *devdelta) {
59
60     /* Get the global ID of the thread. */
61     int tid = get_tid();
62
63     if (tid < numObjs) {
64        int index, i;
65        double dist, min_dist;
66
67        /* find the cluster id that has min distance to object */
68        index = 0;
69
70        min_dist = euclid_dist_2(numCoords, numObjs, numClusters,
71                                 objects, deviceClusters,
72                                 tid, index);
73
74        for (i = 1; i < numClusters; i++) {
75
76           dist = euclid_dist_2(numCoords, numObjs, numClusters,
77                                objects, deviceClusters,
78                                tid, i);
79        /* no need square root */
80        if (dist < min_dist) { /* find the min and its array index */
81           min_dist = dist;
82           index = i;
83        }
84      }
85
86      if (deviceMembership[tid] != index) {
87
88         atomicAdd(devdelta, 1.0);
89      }
90
91      /* assign the deviceMembership to object objectId */
92      deviceMembership[tid] = index;
93   }
94  }
95
96  //
97  // ----------------------------------------
98  //  DATA LAYOUT
99  //
100 //  objects        [numObjs][numCoords]
101 //  clusters       [numClusters][numCoords]
102 //  newClusters    [numClusters][numCoords]
103 //  deviceObjects  [numObjs][numCoords]
104 //  deviceClusters [numClusters][numCoords]
105 //  ----------------------------------------
```

```
106   //
107   /* return an array of cluster centers of size [numClusters][numCoords]        */
108   void kmeans_gpu(double *objects,        /* in: [numObjs][numCoords] */
109                   int numCoords,     /* no. features */
110                   int numObjs,       /* no. objects */
111                   int numClusters,   /* no. clusters */
112                   double threshold,     /* % objects change membership */
113                   long loop_threshold,    /* maximum number of iterations */
114                   int *membership,    /* out: [numObjs] */
115                   double *clusters,    /* out: [numClusters][numCoords] */
116                   int blockSize) {
117     double timing = wtime(), timing_internal, timer_min = 1e42, timer_max = 0;
118     double timing_gpu, timing_cpu, timing_transfers, transfers_time = 0.0, cpu_time = 0.0,
      gpu_time = 0.0;
119     int loop_iterations = 0;
120     int i, j, index, loop = 0;
121     int *newClusterSize; /* [numClusters]: no. objects assigned in each
122                              new cluster */
123     double delta = 0, *dev_delta_ptr;          /* % of objects change their clusters */
124     double **newClusters = (double **) calloc_2d(numClusters, numCoords, sizeof(double));
125
126     double *deviceObjects;
127     double *deviceClusters;
128     int *deviceMembership;
129
130     printf("\n|-------------Naive GPU Kmeans--------------|\n\n");
131
132
133     /* initialize membership[] */
134     for (i = 0; i < numObjs; i++) membership[i] = -1;
135
136     /* need to initialize newClusterSize and newClusters[0] to all 0 */
137     newClusterSize = (int *) calloc(numClusters, sizeof(int));
138     assert(newClusterSize != NULL);
139
140     timing = wtime() - timing;
141     printf("t_alloc: %lf ms\n\n", 1000 * timing);
142     timing = wtime();
143
144     const unsigned int numThreadsPerClusterBlock = (numObjs > blockSize) ? blockSize :
      numObjs;
145     const unsigned int numClusterBlocks = (numObjs + numThreadsPerClusterBlock - 1) /
      numThreadsPerClusterBlock; /* TODO: Calculate Grid size, e.g. number of blocks. */
146     const unsigned int clusterBlockSharedDataSize = 0;
147
148     checkCuda(cudaMalloc(&deviceObjects, numObjs * numCoords * sizeof(double)));
149     checkCuda(cudaMalloc(&deviceClusters, numClusters * numCoords * sizeof(double)));
150     checkCuda(cudaMalloc(&deviceMembership, numObjs * sizeof(int)));
151     checkCuda(cudaMalloc(&dev_delta_ptr, sizeof(double)));
152
153     timing = wtime() - timing;
154     printf("t_alloc_gpu: %lf ms\n\n", 1000 * timing);
155     timing = wtime();
156
```

```
157    checkCuda(cudaMemcpy(deviceObjects, objects,
158                    numObjs * numCoords * sizeof(double), cudaMemcpyHostToDevice));
159    checkCuda(cudaMemcpy(deviceMembership, membership,
160                    numObjs * sizeof(int), cudaMemcpyHostToDevice));
161    timing = wtime() - timing;
162    printf("t_get_gpu: %lf ms\n\n", 1000 * timing);
163    timing = wtime();
164
165    do {
166      timing_internal = wtime();
167
168      /* GPU part: calculate new memberships */
169
170      timing_transfers = wtime();
171
172      checkCuda(cudaMemcpy(deviceClusters, clusters,
173                    numClusters * numCoords * sizeof(double),
174                    cudaMemcpyHostToDevice));
175
176      transfers_time += wtime() - timing_transfers;
177
178      checkCuda(cudaMemset(dev_delta_ptr, 0, sizeof(double)));
179
180      //printf("Launching find_nearest_cluster Kernel with grid_size = %d, block_size = %d,
    shared_mem = %d KB\n", numClusterBlocks, numThreadsPerClusterBlock, clusterBlockSharedDa-
    taSize/1000);
181      timing_gpu = wtime();
182      find_nearest_cluster
183      <<< numClusterBlocks, numThreadsPerClusterBlock, clusterBlockSharedDataSize >>>
184              (numCoords, numObjs, numClusters,
185               deviceObjects, deviceClusters, deviceMembership, dev_delta_ptr);
186
187      cudaDeviceSynchronize();
188      checkLastCudaError();
189      gpu_time += wtime() - timing_gpu;
190      //printf("Kernels complete for itter %d, updating data in CPU\n", loop);
191
192      timing_transfers = wtime();
193
194      checkCuda(cudaMemcpy(membership, deviceMembership,
195                    numObjs * sizeof(int),
196                    cudaMemcpyDeviceToHost));
197
198      checkCuda(cudaMemcpy(&delta, dev_delta_ptr,
199                    sizeof(double),
200                    cudaMemcpyDeviceToHost));
201
202      transfers_time += wtime() - timing_transfers;
203
204      /* CPU part: Update cluster centers*/
205      timing_cpu = wtime();
206      for (i = 0; i < numObjs; i++) {
207        /* find the array index of nestest cluster center */
208        index = membership[i];
```

```
209
210            /* update new cluster centers : sum of objects located within */
211            newClusterSize[index]++;
212            for (j = 0; j < numCoords; j++)
213              newClusters[index][j] += objects[i * numCoords + j];
214          }
215
216          /* average the sum and replace old cluster centers with newClusters */
217          for (i = 0; i < numClusters; i++) {
218            for (j = 0; j < numCoords; j++) {
219              if (newClusterSize[i] > 0)
220                clusters[i * numCoords + j] = newClusters[i][j] / newClusterSize[i];
221              newClusters[i][j] = 0.0;    /* set back to 0 */
222            }
223            newClusterSize[i] = 0;    /* set back to 0 */
224          }
225
226          delta /= numObjs;
227          //printf("delta is %f - ", delta);
228          loop++;
229          //printf("completed loop %d\n", loop);
230          cpu_time += wtime() - timing_cpu;
231
232          timing_internal = wtime() - timing_internal;
233          if (timing_internal < timer_min) timer_min = timing_internal;
234          if (timing_internal > timer_max) timer_max = timing_internal;
235        } while (delta > threshold && loop < loop_threshold);
236
237      timing = wtime() - timing;
238      printf("nloops = %d  : total = %lf ms\n\t-> t_loop_avg = %lf ms\n\t-> t_loop_min = %lf
    ms\n\t-> t_loop_max = %lf ms\n\t"
239             "-> t_cpu_avg = %lf ms\n\t-> t_gpu_avg = %lf ms\n\t-> t_transfers_avg = %lf
    ms\n\n|-------------------------------------------|\n",
240             loop, 1000 * timing, 1000 * timing / loop, 1000 * timer_min, 1000 * timer_max,
241             1000 * cpu_time / loop, 1000 * gpu_time / loop, 1000 * transfers_time / loop);
242
243      char outfile_name[1024] = {0};
244      sprintf(outfile_name, "Execution_logs/silver1-V100_Sz-%lu_Coo-%d_Cl-%d.csv",
245             numObjs * numCoords * sizeof(double) / (1024 * 1024), numCoords, numClusters);
246      FILE *fp = fopen(outfile_name, "a+");
247      if (!fp) error("Filename %s did not open succesfully, no logging performed\n",
    outfile_name);
248      fprintf(fp, "%s,%d,%lf,%lf,%lf\n", "Naive", blockSize, timing / loop, timer_min,
    timer_max);
249      fclose(fp);
250      checkCuda(cudaFree(deviceObjects));
251      checkCuda(cudaFree(deviceClusters));
252      checkCuda(cudaFree(deviceMembership));
253
254      free(newClusters[0]);
255      free(newClusters);
256      free(newClusterSize);
257
258      return;
```

```
259   }
260
261
```