

```

a6/kmeans/kmeans.c

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <mpi.h>
4
5 #include "kmeans.h"
6
7 // square of Euclid distance between two multi-dimensional points
8 inline static double euclid_dist_2(int      numdims, /* no. dimensions */
9                                double * coord1, /* [numdims] */
10                               double * coord2) /* [numdims] */
11 {
12     int i;
13     double ans = 0.0;
14
15     for(i=0; i<numdims; i++)
16         ans += (coord1[i]-coord2[i]) * (coord1[i]-coord2[i]);
17
18     return ans;
19 }
20
21 inline static int find_nearest_cluster(int      numClusters, /* no. clusters */
22                                         int      numCoords, /* no. coordinates */
23                                         double * object, /* [numCoords] */
24                                         double * clusters) /* [numClusters][numCoords] */
25 */
26 {
27     int index, i;
28     double dist, min_dist;
29
30     // find the cluster id that has min distance to object
31     index = 0;
32     min_dist = euclid_dist_2(numCoords, object, clusters);
33
34     for(i=1; i<numClusters; i++) {
35         dist = euclid_dist_2(numCoords, object, &clusters[i*numCoords]);
36         // no need square root
37         if (dist < min_dist) { // find the min and its array index
38             min_dist = dist;
39             index    = i;
40         }
41     }
42     return index;
43 }
44 void kmeans(double * objects,          /* in: [numObjs][numCoords] */
45             int      numCoords,        /* no. coordinates */
46             int      numObjs,          /* no. objects */
47             int      numClusters,       /* no. clusters */
48             double   threshold,        /* minimum fraction of objects that change
membership */
49             long    loop_threshold,    /* maximum number of iterations */
50             int     * membership,       /* out: [numObjs] */

```

```

51         double * clusters)          /* out: [numClusters][numCoords] */
52 {
53     int i, j;
54     int index, loop=0;
55     double timing = 0;
56
57     /* Every variable has its "rank_" version, which is used to store local data,
58      * and its "new" version, which is used to store global data.
59      */
60     double rank_delta, delta = 0;           // fraction of objects whose clusters
change in each loop
61     int * rank_newClusterSize, * newClusterSize; // [numClusters]: no. objects assigned in
each new cluster
62     double * rank_newClusters, *newClusters;    // [numClusters][numCoords]
63
64     // Get rank of this process
65     int rank;
66     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
67
68     // initialize membership
69     for (i=0; i<numObjs; i++)
70         membership[i] = -1;
71
72     // initialize rank_newClusterSize and rank_newClusters to all 0
73     rank_newClusterSize = (typeof(rank_newClusterSize)) malloc(numClusters,
sizeof(*rank_newClusterSize));
74     rank_newClusters    = (typeof(rank_newClusters)) malloc(numClusters * numCoords,
sizeof(*rank_newClusters));
75     newClusterSize      = (typeof(newClusterSize)) malloc(numClusters,
sizeof(*newClusterSize));
76     newClusters         = (typeof(newClusters)) malloc(numClusters * numCoords,
sizeof(*newClusters));
77
78     timing = wtime();
79     do {
80         // before each loop, set cluster data to 0
81         for (i=0; i<numClusters; i++) {
82             for (j=0; j<numCoords; j++)
83                 rank_newClusters[i*numCoords + j] = 0.0;
84             rank_newClusterSize[i] = 0;
85         }
86
87         rank_delta = 0.0;
88
89         for (i=0; i<numObjs; i++) {
90             // find the array index of nearest cluster center
91             index = find_nearest_cluster(numClusters, numCoords, &objects[i*numCoords],
clusters);
92
93             // if membership changes, increase rank_delta by 1
94             if (membership[i] != index)
95                 rank_delta += 1.0;
96
97             // assign the membership to object i

```

```

98         membership[i] = index;
99
100        // update new cluster centers : sum of objects located within
101        rank_newClusterSize[index]++;
102        for (j=0; j<numCoords; j++)
103            rank_newClusters[index*numCoords + j] += objects[i*numCoords + j];
104    }
105
106    /*
107     * TODO: Perform reduction of cluster data (rank_newClusters, rank_newClusterSize)
108     * from local arrays to shared.
109     */
110    MPI_Allreduce(rank_newClusters, newClusters, numClusters * numCoords, MPI_DOUBLE,
111    MPI_SUM, MPI_COMM_WORLD);
112    MPI_Allreduce(rank_newClusterSize, newClusterSize, numClusters, MPI_INT, MPI_SUM,
113    MPI_COMM_WORLD);
114
115
116    // average the sum and replace old cluster centers with newClusters
117    for (i=0; i<numClusters; i++) {
118        if (newClusterSize[i] > 0) {
119            for (j=0; j<numCoords; j++) {
120                clusters[i*numCoords + j] = newClusters[i*numCoords + j] /
121                newClusterSize[i];
122            }
123        }
124    }
125
126    /*
127     * TODO: Perform reduction from rank_delta variable to delta variable, that will
128     * be used for convergence check.
129     */
130    MPI_Allreduce(&rank_delta, &delta, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
131
132
133    // Get fraction of objects whose membership changed during this loop. This is used
134    // as a convergence criterion.
135    delta /= numObjs;
136
137    loop++;
138    //printf("\r\tcompleted loop %d", loop);
139    //fflush(stdout);
140
141    } while (delta > threshold && loop < loop_threshold);
142
143    timing = wtime() - timing;
144    if (rank == 0) fprintf(stdout, "      nloops = %3d    (total = %7.4fs)  (per loop =
145    %7.4fs)\n", loop, timing, timing/loop);
146
147
148    free(rank_newClusters);
149    free(rank_newClusterSize);
150    free(newClusters);
151    free(newClusterSize);
152}
153

```