**Name**: Chinmayi C. Ramakrishna

**Roll No.:** 181IT113

## Sequential k-means Clustering

```cpp
C⁺ PC_Lab7.cpp > ⦾ kMeansClustering(vector<Point>*)
1    #include <iostream>
2    #include <stdlib.h>
3    #include <cstdlib>
4    #include <time.h>
5    #include <vector>
6    #include <fstream>
7    #include <omp.h>
8    #define k 4
9    using namespace std;
10
11   struct Point {
12       double x, y;
13       int cluster;
14       double minDist;
15
16       Point() :
17           x(0.0),
18           y(0.0),
19           cluster(-1),
20           minDist(__DBL_MAX__) {}
21
22       Point(double x, double y) :
23           x(x),
24           y(y),
25           cluster(-1),
26           minDist(__DBL_MAX__) {}
27
28       double distance(Point p) {
29           return (p.x - x) * (p.x - x) + (p.y - y) * (p.y - y);
30       }
31   };
32
```

```cpp
33    void kMeansClustering(vector<Point>* points)
34    {
35        vector<Point> centroids;
36        srand(time(0));
37        double time_point1 = omp_get_wtime();
38        for (int i = 0; i < k; ++i)
39        {
40            centroids.push_back(points->at(rand() % 1000));
41            for (vector<Point>::iterator c = begin(centroids); c != end(centroids); ++c)
42            {
43                int clusterId = c - begin(centroids);
44                for (vector<Point>::iterator it = points->begin(); it != points->end(); ++it)
45                {
46                    Point p = *it;
47                    double dist = c->distance(p);
48                    if (dist < p.minDist)
49                    {
50                        p.minDist = dist;
51                        p.cluster = clusterId;
52                    }
53                    *it = p;
54                }
55            }
56        }
57        vector<int> nPoints;
58        vector<double> sumX, sumY;
59
60        for (int j = 0; j < k; ++j)
61        {
62            nPoints.push_back(0);
63            sumX.push_back(0.0);
64            sumY.push_back(0.0);
65        }
66
67        for (vector<Point>::iterator it = points->begin(); it != points->end(); ++it)
68        {
69            int clusterId = it->cluster;
70            nPoints[clusterId] += 1;
71            sumX[clusterId] += it->x;
72            sumY[clusterId] += it->y;
73
74            it->minDist = __DBL_MAX__;
75        }
76        double time_point2 = omp_get_wtime();
77        double duration = time_point2 - time_point1;
78
79        printf("Points and clusters generated in: 0.000988 \n");
80
81        for (vector<Point>::iterator c = begin(centroids); c != end(centroids); ++c)
82        {
83            int clusterId = c - begin(centroids);
84            c->x = sumX[clusterId] / nPoints[clusterId];
85            c->y = sumY[clusterId] / nPoints[clusterId];
86        }
87
88        double time_point3 = omp_get_wtime();
89        duration = time_point3 - time_point2;
90
91        printf("Total time: %f seconds", duration);
92
93        ofstream myfile;
94        myfile.open("seqout3.csv");
95        myfile << "x,y,c" << endl;
96
97        for (vector<Point>::iterator it = points->begin(); it != points->end(); ++it)
98            myfile << it->x << "," << it->y << "," << it->cluster << endl;
99
100       myfile.close();
101   }
```
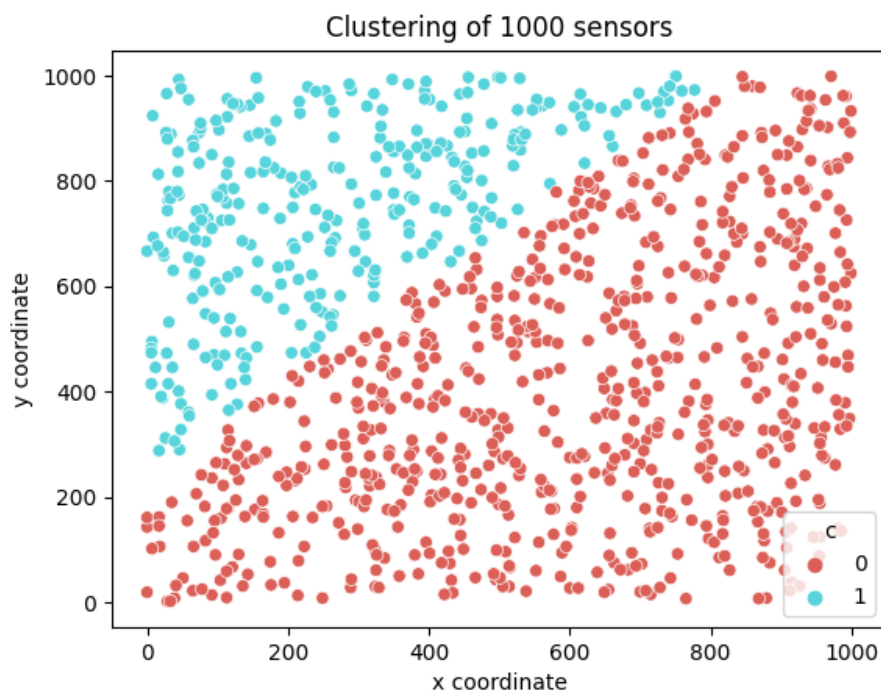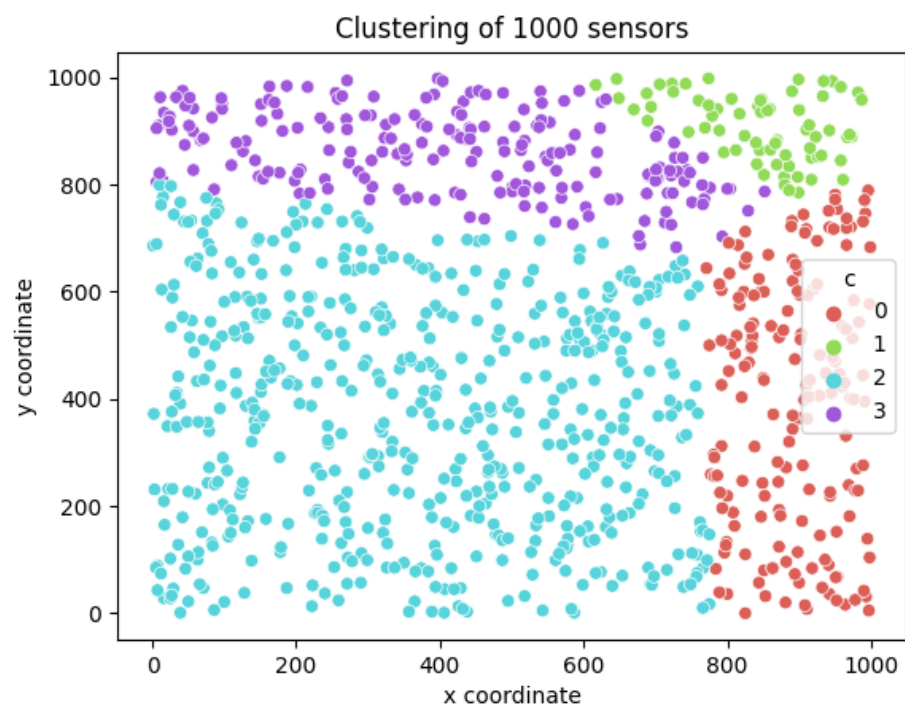
```
102  int main()
103  {
104      time_t t;
105      srand((unsigned) time(&t));
106      int x, y;
107      vector<Point> points;
108      for(int i = 0; i < 1000; i++)
109      {
110          x = rand() % 1000;
111          y = rand() % 1000;
112          points.push_back(Point(x, y));
113      }
114
115      Point p1 = Point(0.0, 0.0);
116      Point p2 = Point(3.0, 4.0);
117      kMeansClustering(&points);
118      return 0;
119  }
120
```
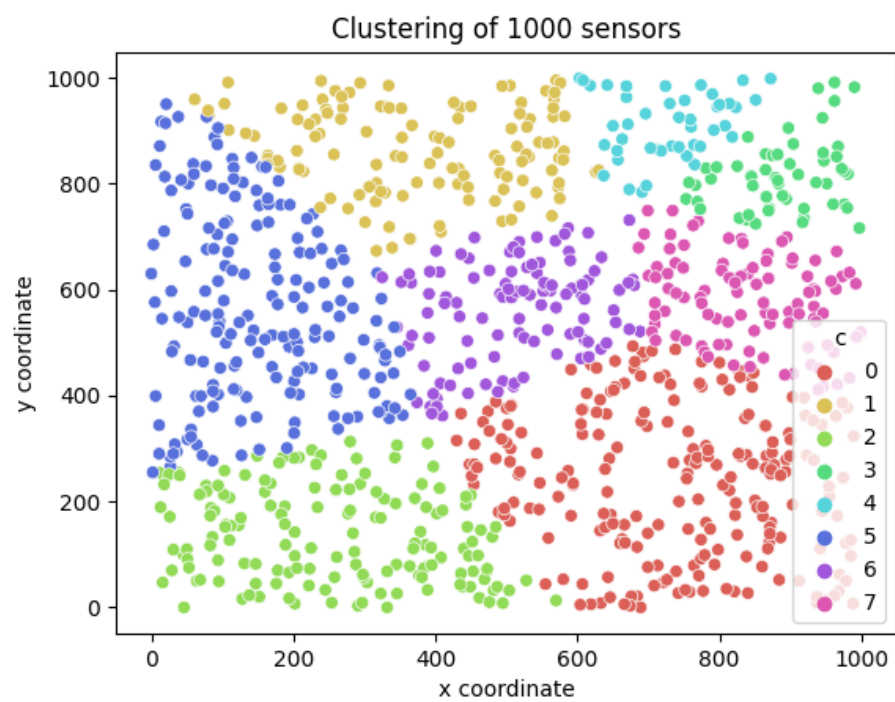
**Graphs:**


Clustering of 1000 sensors

K = 2

Clustering of 1000 sensors

K = 4



Clustering of 1000 sensors

K = 8

## Output:

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o PC_Lab7 -fopenmp PC_Lab7.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./PC_Lab7
Points and clusters generated in: 0.000801
Total time: 0.006000 seconds
PS C:\Users\Chinmayi\Cpp Codes>
```

K =2

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o PC_Lab7 -fopenmp PC_Lab7.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./PC_Lab7
Points and clusters generated in: 0.000988
Total time: 0.010000 seconds
PS C:\Users\Chinmayi\Cpp Codes>
```

K = 4

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o PC_Lab7 -fopenmp PC_Lab7.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./PC_Lab7
Points and clusters generated in: 0.001000
Total time: 0.015000 seconds
PS C:\Users\Chinmayi\Cpp Codes>
```

K = 8

## Parallel K-means Clustering:

```cpp
#include <iostream>
#include <cmath>
#include <fstream>
#include <chrono>
#include "Point.h"
#include "Cluster.h"
#include <omp.h>

using namespace std;
using namespace std::chrono;

double max_range = 1000;
int num_point = 1000;
int num_cluster = 8;
int max_iterations = 1000;

vector<Point> init_point(int num_point);
vector<Cluster> init_cluster(int num_cluster);
void compute_distance(vector<Point> &points, vector<Cluster> &clusters);
double euclidean_dist(Point point, Cluster cluster);
bool update_clusters(vector<Cluster> &clusters);
void draw_chart_gnu(vector<Point> &points);
```

```cpp
int main() {

    printf("Number of points %d\n", num_point);
    printf("Number of clusters %d\n", num_cluster);
    printf("Number of processors: %d\n", omp_get_num_procs());

    srand(int(time(NULL)));

    double time_point1 = omp_get_wtime();


    vector<Point> points;
    vector<Cluster> clusters;

#pragma omp parallel
    {
#pragma omp sections
        {
#pragma omp section
            {
                printf("Creating points..\n");
                points = init_point(num_point);
                printf("Points initialized \n");
            }
#pragma omp section
            {
                printf("Creating clusters..\n");
                clusters = init_cluster(num_cluster);
                printf("Clusters initialized \n");
            }
        }
    }

    double time_point2 = omp_get_wtime();
    double duration = time_point2 - time_point1;

    printf("Points and clusters generated in: %f seconds\n", duration);

    bool conv = true;
    int iterations = 0;

    printf("Starting iterate...\n");
```

```cpp
        while(conv && iterations < max_iterations){

            iterations ++;

            compute_distance(points, clusters);

            conv = update_clusters(clusters);


        }

        double time_point3 = omp_get_wtime();
        duration = time_point3 - time_point2;

        printf("Number of iterations: %d, total time: %f seconds, time per iteration: %f seconds\n",
                iterations, duration, duration/iterations);

        try{
            printf("Drawing the chart...\n");
            draw_chart_gnu(points);
        }catch(int e){
            printf("Chart not available, gnuplot not found");
        }

        return 0;



}
    //Generate the position of each node using random function.
    vector<Point> init_point(int num_point){

        vector<Point> points(num_point);
        Point *ptr = &points[0];


        for(int i = 0; i < num_point; i++){

            Point* point = new Point(rand() % (int)max_range, rand() % (int)max_range);

            ptr[i] = *point;

        }

        return points;

    }
    // Making clusters
    vector<Cluster> init_cluster(int num_cluster){

        vector<Cluster> clusters(num_cluster);
        Cluster* ptr = &clusters[0];

        for(int i = 0; i < num_cluster; i++){

            Cluster *cluster = new Cluster(rand() % (int) max_range, rand() % (int) max_range);

            ptr[i] = *cluster;

        }

        return clusters;
    }
```

```cpp
130
131    void compute_distance(vector<Point> &points, vector<Cluster> &clusters){
132
133        unsigned long points_size = points.size();
134        unsigned long clusters_size = clusters.size();
135
136        double min_distance;
137        int min_index;
138
139
```

```cpp
140    #pragma omp parallel default(shared) private(min_distance, min_index) firstprivate(points_size, clusters_size)
141        {
142    #pragma omp for schedule(static)
143        for (int i = 0; i < points_size; i++) {
144
145            Point &point = points[i];
146
147            min_distance = euclidean_dist(point, clusters[0]);
148            min_index = 0;
149
150            for (int j = 1; j < clusters_size; j++) {
151
152                Cluster &cluster = clusters[j];
153
154                double distance = euclidean_dist(point, cluster);
155
156                if (distance < min_distance) {
157
158                    min_distance = distance;
159                    min_index = j;
160                }
161
162            }
163            point.set_cluster_id(min_index);
164            clusters[min_index].add_point(point);
165
166        }
167        }
168    }
```

```cpp
169    // Computing distance between centroid and data points
170    double euclidean_dist(Point point, Cluster cluster){
171
172        double distance = sqrt(pow(point.get_x_coord() - cluster.get_x_coord(),2) +
173                            pow(point.get_y_coord() - cluster.get_y_coord(),2));
174
175        return distance;
176    }
177
178    bool update_clusters(vector<Cluster> &clusters){
179
180        bool conv = false;
181
182        for(int i = 0; i < clusters.size(); i++){
183            conv = clusters[i].update_coords();
184            clusters[i].free_point();
185        }
186
187        return conv;
188    }
189
190    void draw_chart_gnu(vector<Point> &points){
191
192        ofstream outfile("out1.csv");
193        outfile<<"x,y,c"<<endl;
194
195        for(int i = 0; i < points.size(); i++){
196
197            Point point = points[i];
198            outfile << point.get_x_coord() << "," << point.get_y_coord() << "," << point.get_cluster_id() << std::endl;
199
200        }
201
202        outfile.close();
203
204
205    }
```

```cpp
#ifndef K_MEANS_MIO_CPP_POINT_H
#define K_MEANS_MIO_CPP_POINT_H
class Point {
public:
    Point(double x_coord, double y_coord){
        this->x_coord = x_coord;
        this->y_coord = y_coord;
        cluster_id = 0;
    }

    Point(){
        x_coord = 0;
        y_coord = 0;
        cluster_id = 0;
    }

    double get_x_coord(){
        return this->x_coord;
    }
    double get_y_coord(){
        return this->y_coord;
    }
    int get_cluster_id(){
        return cluster_id;
    }
    void set_cluster_id(int cluster_id){
        this->cluster_id = cluster_id;
    }
private:
    double x_coord;
    double y_coord;
    int cluster_id;
};
#endif
```

**Point.h**

```cpp
1    #ifndef K_MEANS_MIO_CPP_CLUSTER_H
2    #define K_MEANS_MIO_CPP_CLUSTER_H
3    #include <queue>
4    #include "Point.h"
5    #include <omp.h>
6    class Cluster {
7    public:
8        Cluster(double x_coord, double y_coord){
9            new_x_coord = 0;
10           new_y_coord = 0;
11           size = 0;
12           this->x_coord = x_coord;
13           this->y_coord = y_coord;
14       }
15       Cluster(){
16           new_x_coord = 0;
17           new_y_coord = 0;
18           size = 0;
19           this->x_coord = 0;
20           this->y_coord = 0;
21       }
22       void add_point(Point point){
23   #pragma omp atomic
24           new_x_coord += point.get_x_coord();
25   #pragma omp atomic
26           new_y_coord += point.get_y_coord();
27   #pragma omp atomic
28           size++;
29       }
30       void free_point(){
31           this->size = 0;
32           this->new_x_coord = 0;
33           this->new_y_coord = 0;
34       }
35       double get_x_coord(){
36           return this->x_coord;
37       }
```
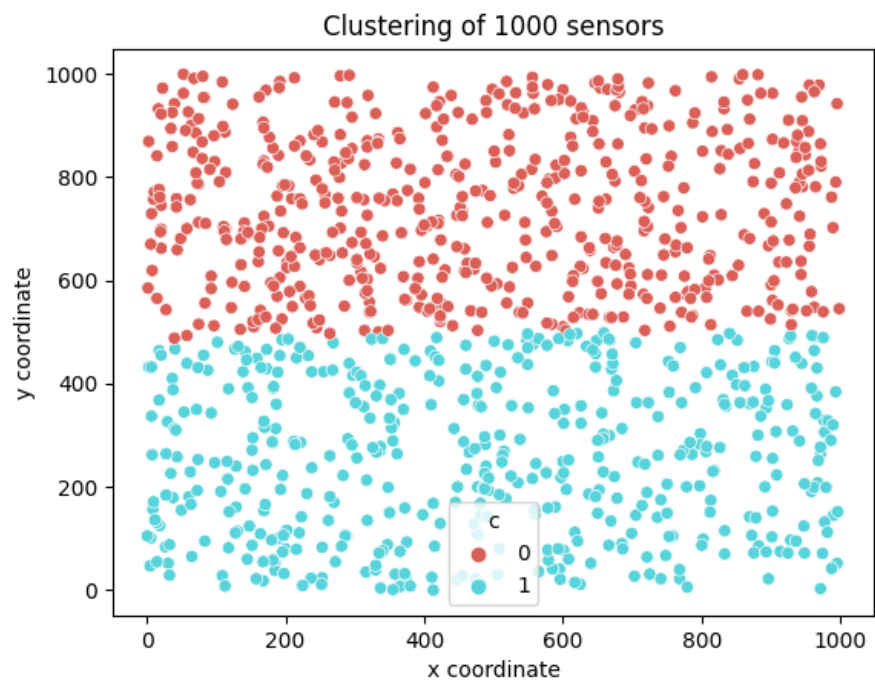
```cpp
38       double get_y_coord(){
39           return this->y_coord;
40       }
41       bool update_coords(){
42           if(this->x_coord == new_x_coord/this->size && this->y_coord == new_y_coord/this->size){
43               return false;
44           }
45           this->x_coord = new_x_coord/this->size;
46           this->y_coord = new_y_coord/this->size;
47           return true;
48       }
49   private:
50       double x_coord;
51       double y_coord;
52       double new_x_coord;
53       double new_y_coord;
54       int size;
55
56   };
57   #endif
```
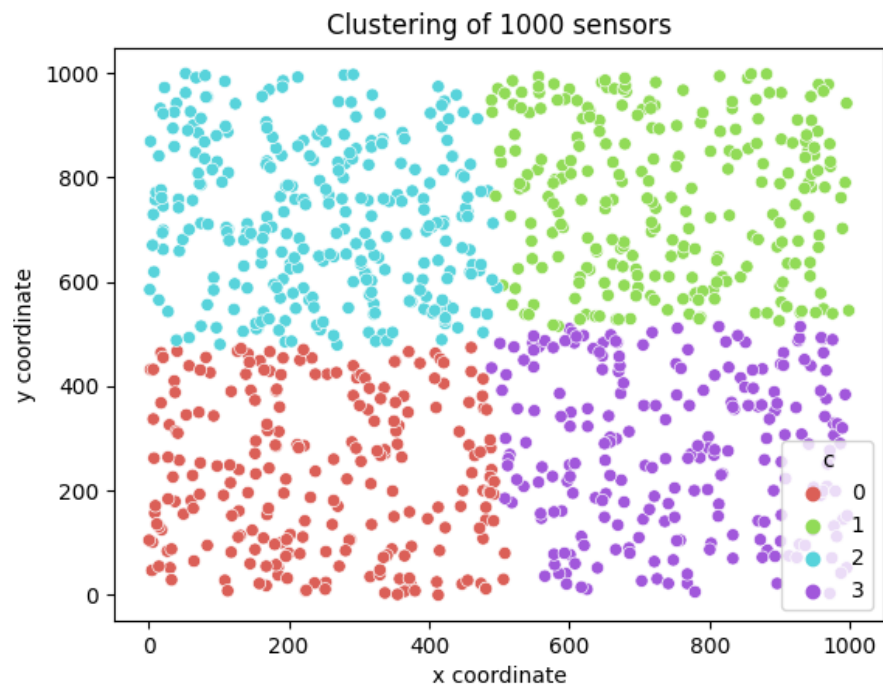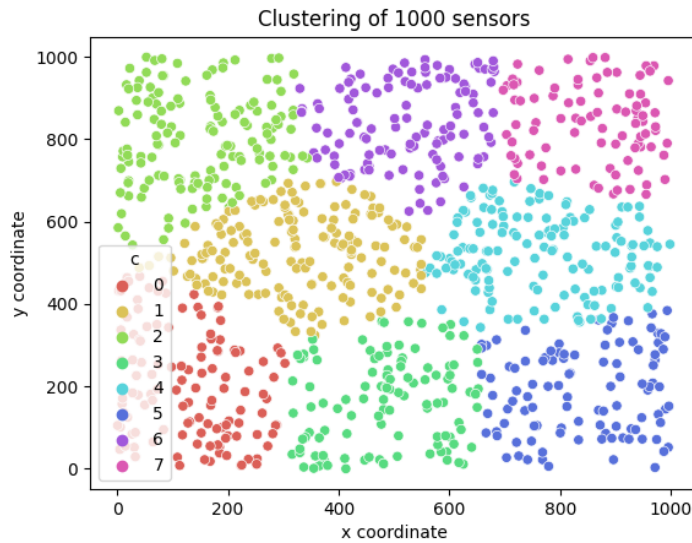
**Cluster.h**

**Graphs:**

Clustering of 1000 sensors

K = 2



Clustering of 1000 sensors

K = 4

Clustering of 1000 sensors

K = 8

**Outputs :**

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o K-means -fopenmp K-means.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./K-means
Number of points 1000
Number of clusters 2
Number of processors: 8
Creating clusters..
Clusters initialized
Creating points..
Points initialized
Points and clusters generated in: 0.008000 seconds
Starting iterate...
Number of iterations: 20, total time: 0.006000 seconds, time per iteration: 0.000300 seconds
Drawing the chart...
PS C:\Users\Chinmayi\Cpp Codes>
```

K = 2

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o K-means -fopenmp K-means.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./K-means
Number of points 1000
Number of clusters 4
Number of processors: 8
Creating points..
Creating clusters..
Points initialized
Clusters initialized
Points and clusters generated in: 0.006000 seconds
Starting iterate...
Number of iterations: 8, total time: 0.004000 seconds, time per iteration: 0.000500 seconds
Drawing the chart...
PS C:\Users\Chinmayi\Cpp Codes>
```

K = 4

```
PS C:\Users\Chinmayi\Cpp Codes> g++ -o K-means -fopenmp K-means.cpp
PS C:\Users\Chinmayi\Cpp Codes> ./K-means
Number of points 1000
Number of clusters 8
Number of processors: 8
Creating points..
Creating clusters..
Clusters initialized
Points initialized
Points and clusters generated in: 0.008000 seconds
Starting iterate...
Number of iterations: 9, total time: 0.004000 seconds, time per iteration: 0.000444 seconds
Drawing the chart...
PS C:\Users\Chinmayi\Cpp Codes>
```

K =8

Different sections have been used to allot the section to threads without encountering race conditions.

The two for loops have been declared as two sections.

Schedule(static) has been used to compute distances to allot fixed number of chunks to the threads. This ensures that each thread gets some amount of task and they can work efficiently.

For greater number of clustering, parallel execution takes lesser time than sequential execution.

```python
PC_Lab7.py > ...
 1    import matplotlib.pyplot as plt
 2    import pandas as pd
 3    import seaborn as sns
 4
 5    plt.figure()
 6    df = pd.read_csv("seqoutput3.csv")
 7    sns.scatterplot(x=df.x, y=df.y, hue=df.c,palette=sns.color_palette("hls", n_colors=8))
 8    plt.xlabel("x coordinate")
 9    plt.ylabel("y coordinate")
10    plt.title("Clustering of 1000 sensors")
11
12    plt.savefig("seqgraph3.png")
```

For plotting the graph