

Clustering & K-means Algorithm

ABSTRACT

Clustering is one of the most important tools in data mining. It can cluster data record without any prior information. K-means is one of the clustering algorithm, which is popular simple and effective. I will implement the K-means algorithm with Matlab code and compare it with the Matlab's statistic tools' K-means algorithm and try to show the different in effectiveness and efficiency aspects. The data I will use to test algorithm from Professor Zhao's website(<http://sse.tongji.edu.cn/zhaoginpei/Datasets/>).

Keywords: Data Mining, Clustering, K-means Algorithm, Centroids, Effectiveness, Efficiency, Local optimal, Voronoi.

INTRODUCTION

Cluster & Clustering:

Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters). It is a main task of exploratory data mining, and a common technique for statistical data analysis, used in many fields, including machine learning, pattern recognition, image analysis, information retrieval, and bioinformatics.

K-means Algorithm:

k-means clustering is a method of vector quantization, originally from signal processing, that is popular for cluster analysis in data mining. k-means clustering aims to partition n observations into k clusters in which each observation belongs to the cluster with the

nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells.

Given a set of observations (x_1, x_2, \dots, x_n) , where each observation is a d -dimensional real vector, k -means clustering aims to partition the n observations into k ($\leq n$) sets $S = \{S_1, S_2, \dots, S_k\}$ so as to minimize the within-cluster sum of squares (WCSS) (sum of distance functions of each point in the cluster to the K center). In other words, its objective is to find: where μ_i is the mean of points in S_i .

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2$$

Figure 1.1

ALGORITHM & IMPLEMENT

2.1 K-means algorithm:

Given an initial set of k means $m_1(1), \dots, m_k(1)$ (see below), the algorithm proceeds by alternating between two steps:

Assignment step: Assign each observation to the cluster whose mean yields the least within-cluster sum of squares (WCSS). Since the sum of squares is the squared Euclidean distance, this is intuitively the "nearest" mean.[8] (Mathematically, this means partitioning the observations according to the Voronoi diagram generated by the means).

$$S_i^{(t)} = \{x_p : \|x_p - m_i^{(t)}\|^2 \leq \|x_p - m_j^{(t)}\|^2 \forall j, 1 \leq j \leq k\},$$

where each x_p is assigned to exactly one $S_i^{(t)}$, even if it could be assigned to two or more of them.

Figure 2.1

Update step: Calculate the new means to be the centroids of the observations in the new clusters.

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

Figure 2.2

Since the arithmetic mean is a least-squares estimator, this also minimizes the within-cluster sum of squares (WCSS) objective.

The algorithm has converged when the assignments no longer change. Since both steps optimize the WCSS objective, and there only exists a finite number of such partitionings, the algorithm must converge to a (local) optimum. There is no guarantee that the global optimum is found using this algorithm.

2.2 Matlab implementation code:

```

1  function [cluster,centroid] = KmeansAlgriothm(k,data)
2  -   len = length(data);
3  -   tmp=randperm(len);
4  -   index = tmp(1:k);
5  -   centroidX = data(index,1);
6  -   centroidY = data(index,2);
7  -   cluster = zeros(len,1);
8  -   % 随机选取k个点作为初始值
9  -   centroid = [centroidX,centroidY];
10 -   while 1
11 -       % 保留原来的centroid
12 -       old = centroid;
13 -       % 计算每个点的distance
14 -       for i = 1:len;
15 -           distance = zeros(k,1);
16 -           for j = 1:k
17 -               distance(j) = norm(data(i,:)-centroid(j,:));
18 -           end
19 -           [~,ind] = min(distance);
20 -           cluster(i,1) = ind;
21 -       end
22 -       % 更新centroid
23 -       for m =1:k
24 -           centroid(m,:)= mean(data(cluster==m,:));
25 -       end
26 -       if (isequal(old,centroid))
27 -           break;
28 -       end
29 -   end
30 - end

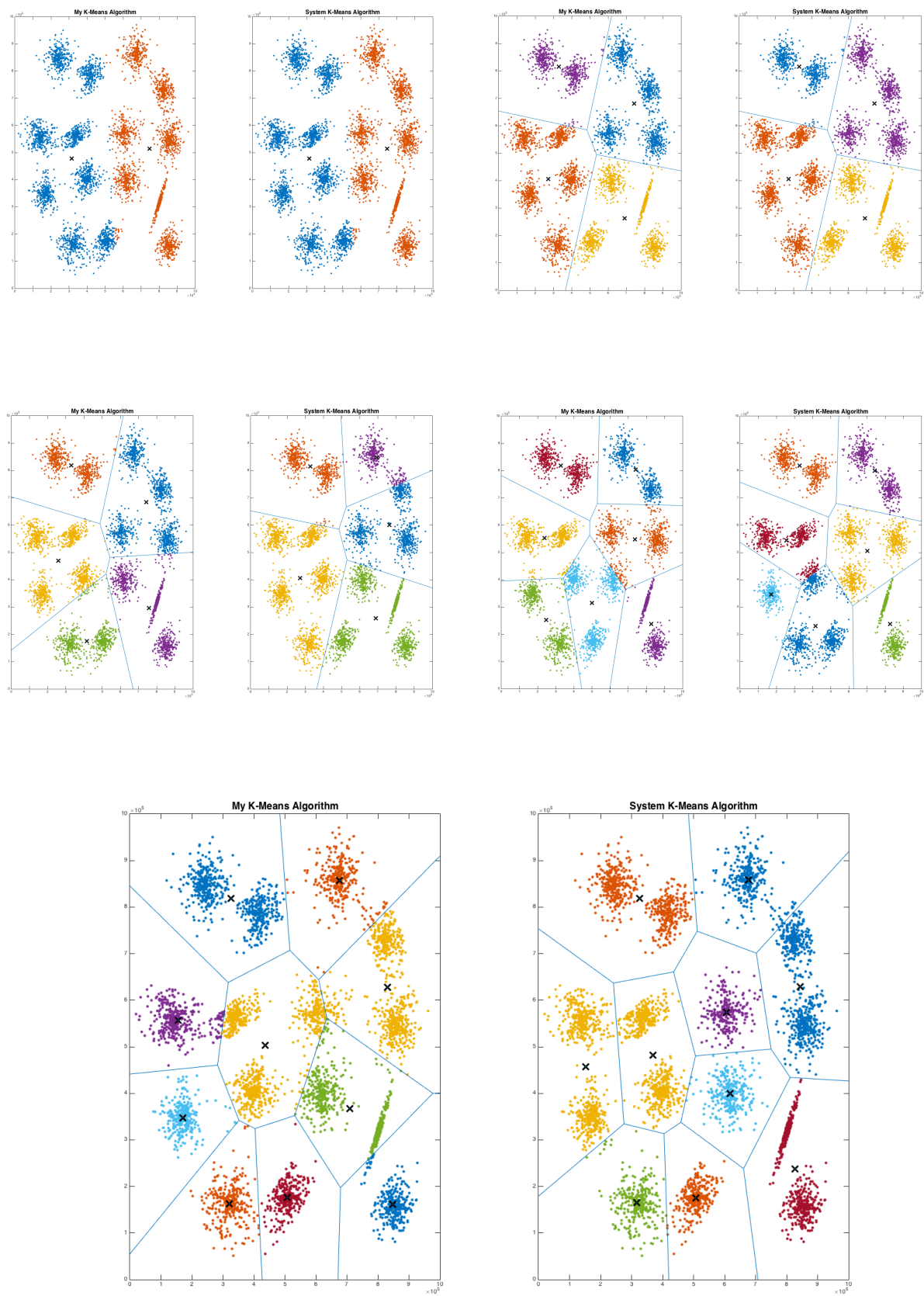
```

EXPERIEMENT

3.1 My K-means Algorithm Result & Matlab Statistic Tools' K-means Algorithm Result

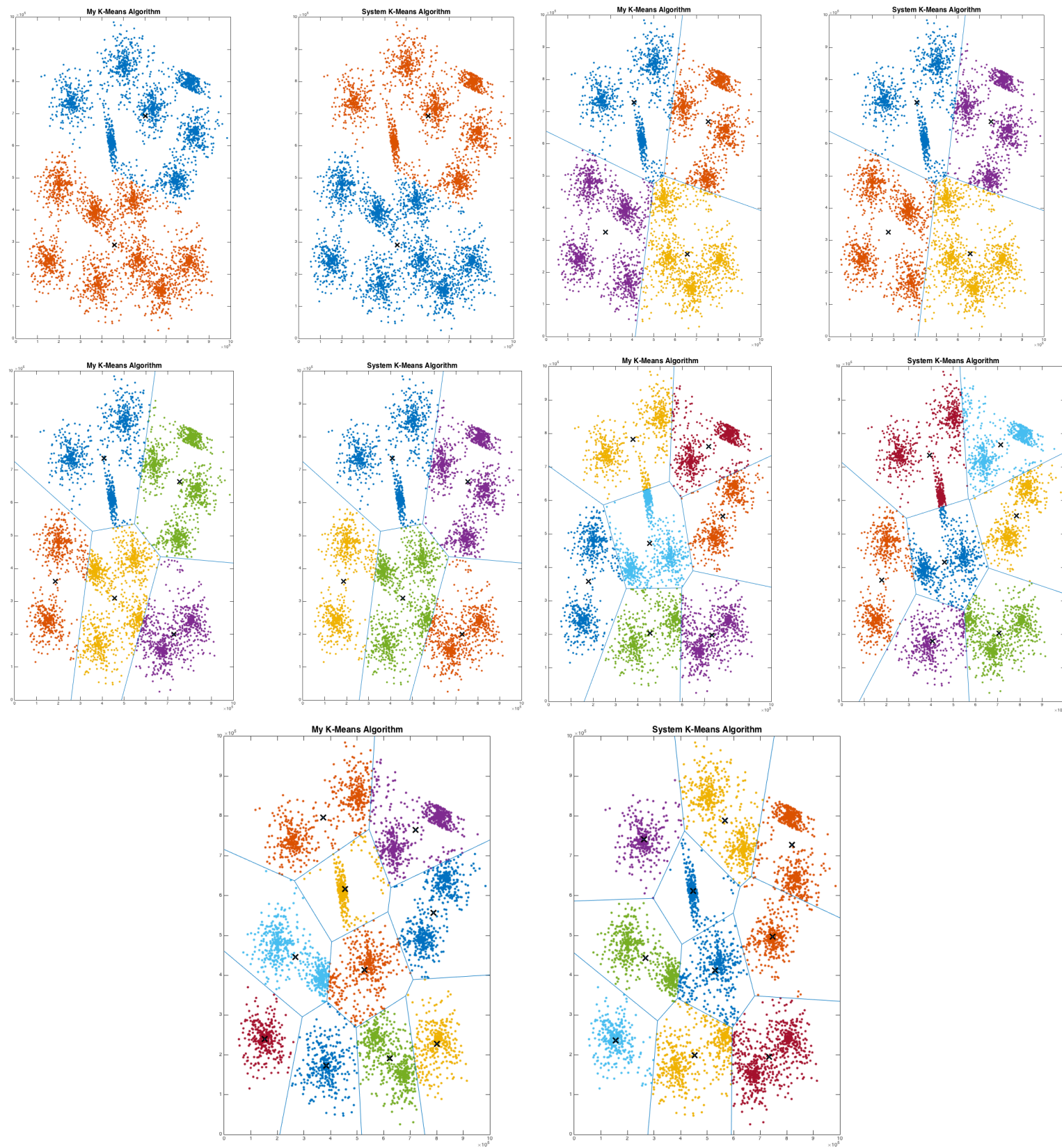
Data Set 1:

k=2 4 5 7 10



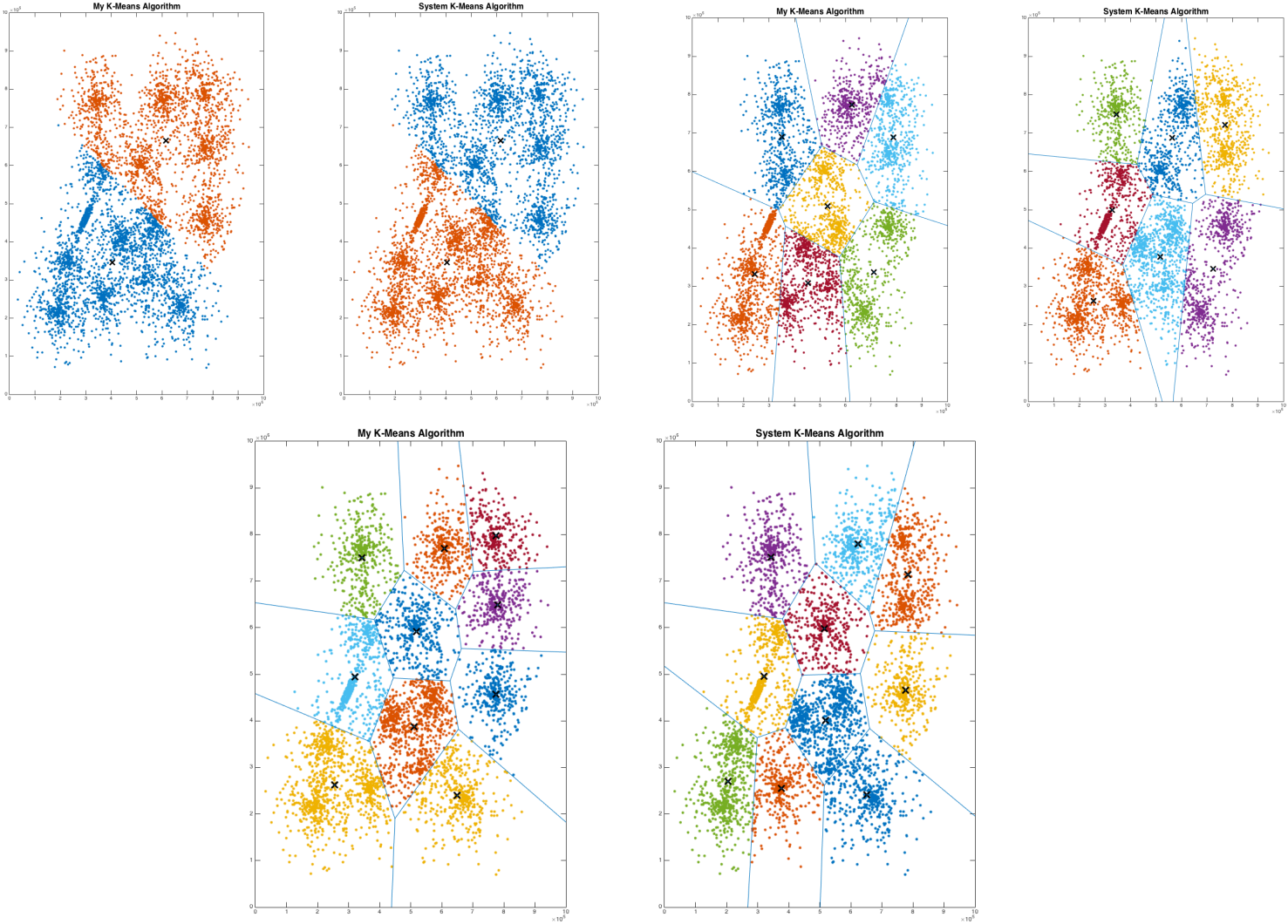
Data Set 2:

k=2 4 5 7 10



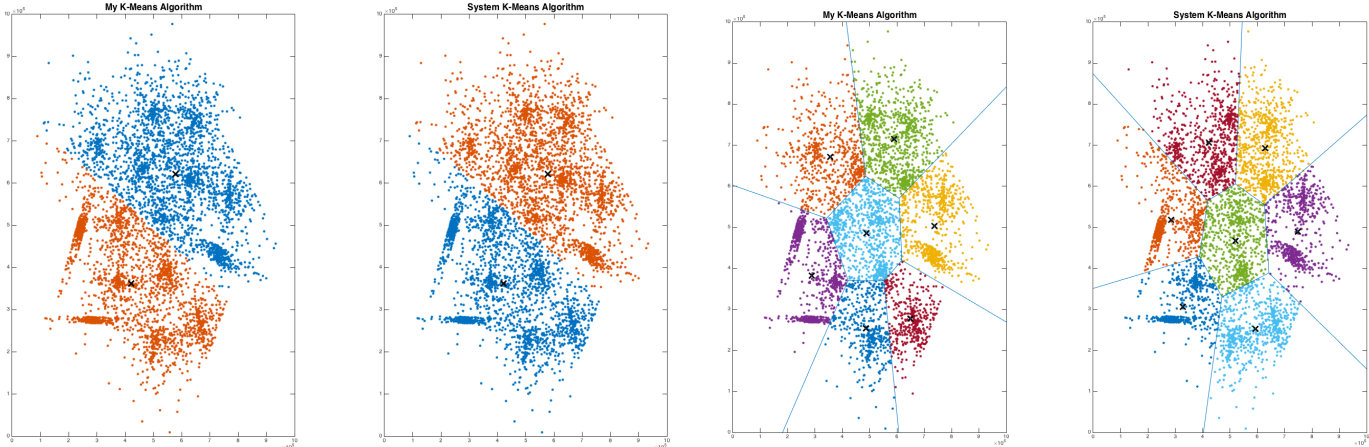
Data Set 3:

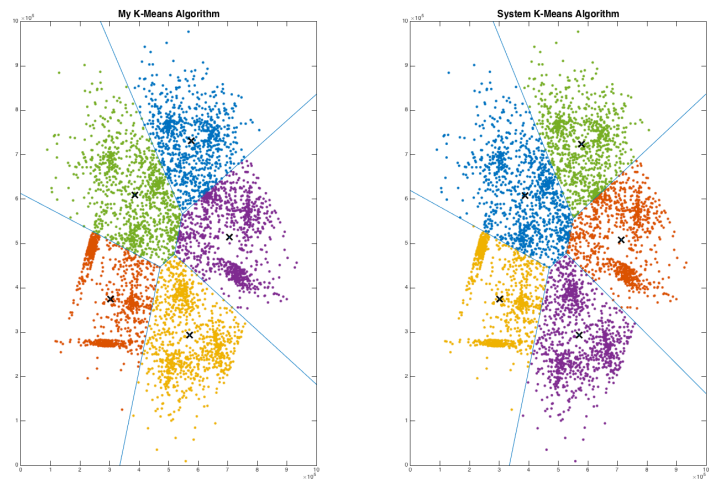
k=2 7 10



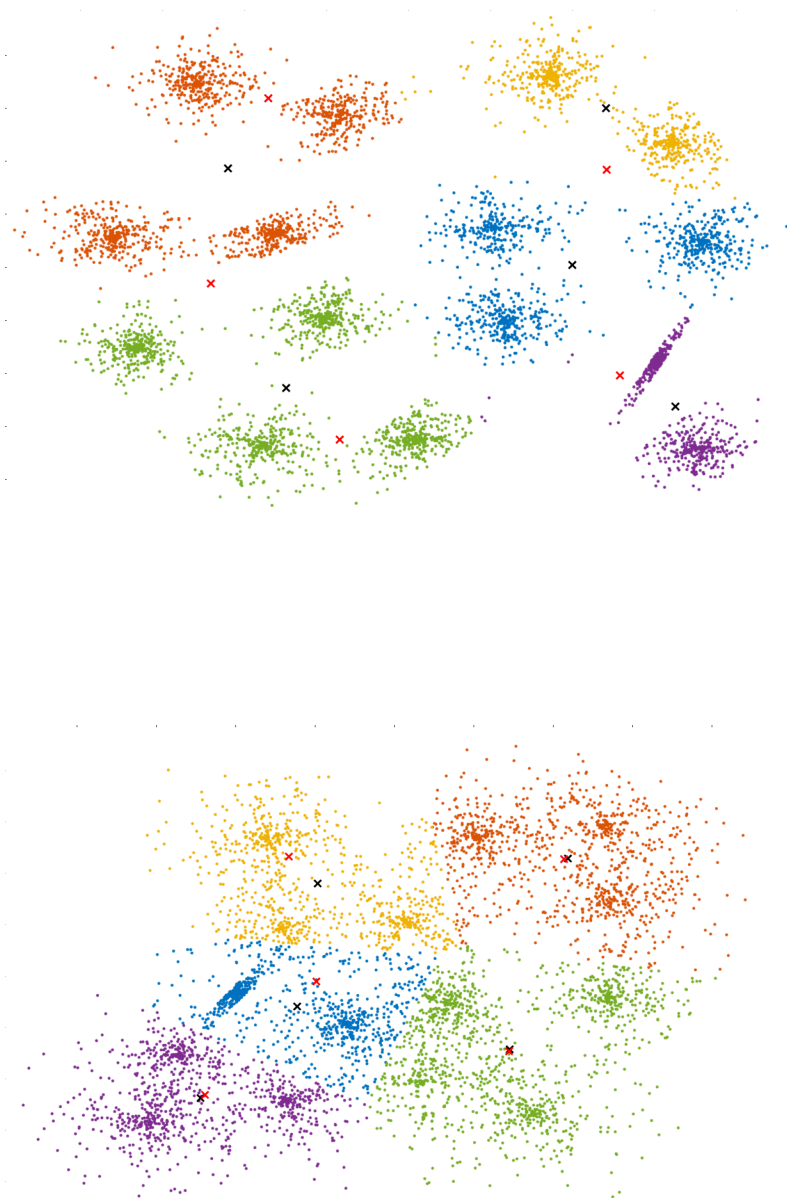
Data Set 4:

k=2 7 5

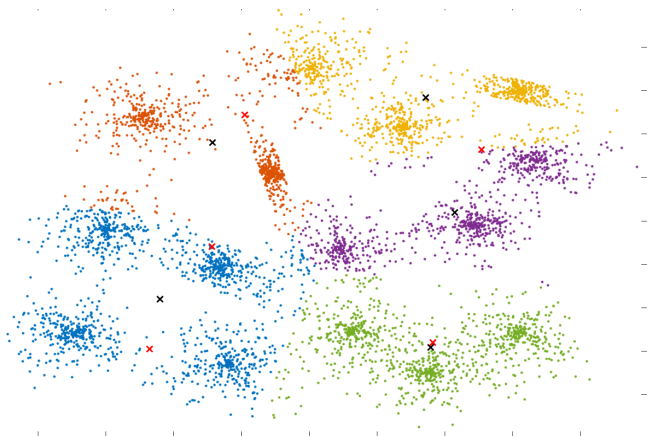


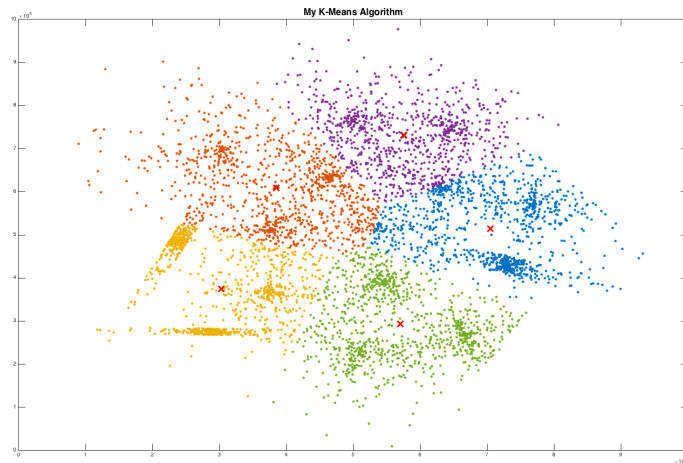


3.2 Compare (k = 5)



black centroids—my K-means
red centroids — statistic Tools K-means





CONCLUSION

1. K-means algorithm will finally reach a local optimal result but this local optimal may not be the global optimal result.
2. K-means algorithm is a NP-hard problem with time consuming.
3. According to random initial centroids, it will show different result. Due to the fact that there may not be one local optimal result.
4. Compare my K-means with statistical tools K-means algorithm.

- **Effectiveness:**

- We all reach one local optimal result finally, however, depend on the initial centroids, the result may be just the same or have some difference (the difference may be very huge, but both of the two results are the local optimal result of this data). In short, it is easy to get the same effectiveness with the system's K-means algorithm as long as you write the right code.

- **Efficiency:**

- Even though I try hard to improve my code's performance in the efficiency, I still have to admit that the system's algorithm performs better than mine. I try to avoid the loop, use the system's calculator method instead of writing by myself, however, it still can't get faster than the system's algorithm. When K's value is less than 10, it is hard to find that my algorithm is slower than the system's. But, with the increase of K's value, the system's algorithm speed still can keep fast, and my algorithm's performance can be that good.

REFERENCE:

- [1] Wikipedia, Cluster analysis, https://en.wikipedia.org/wiki/Cluster_analysis
- [2] Professor ZHAO's homepage, <http://sse.tongji.edu.cn/zhaoqinpei/Datasets/>
- [3] Matlab's official document about K-means, <http://cn.mathworks.com/help/stats/kmeans.html>