

APPENDIX A

Algorithm 1 is shown below:

Algorithm 1: The GB generation algorithm

```

Input: Dataset  $X = \{x_1, x_2, \dots, x_n\}$ , bandwidth  $h$ ,  

merging threshold  $\varepsilon$   

Output: GB Set  $GBs = \{gb_1, gb_2, \dots, gb_m\}$ 
1 Initialize:  $n = |X|, gb = X, GBs = \emptyset, M = \emptyset, P = \emptyset;$   

2 Create an empty queue  $Q$ , add  $gb$  to  $Q$ ;  

3 while  $Q$  is not empty do  

4   Take the first element  $gb$  from  $Q$  and remove it from  $Q$ ;  

5   if  $|gb| > \sqrt{n}$  then  

6     for each  $o_j \in gb$  do  

7       Calculate: kernel density estimate  $\hat{f}(x_j)$  using  

Eq. (6);  

8       Calculate: covariance using Eq. (16);  

9       Calculate: effective dimension  $D_{eff}$  using Eq.  

(15);  

10      Calculate: adaptive bandwidth  $h_i$  using Eq.  

(17);  

11      Calculate: mean shift vector  $m_{h_i}(o_j)$  using Eq.  

(7);  

12      Update: sample position  


$$o_j^{(t+1)} = o_j^{(t)} + m_{h_i}(o_j^{(t)})$$
 using Eq. (12);  

13      if  $\|o_j^{(t+1)} - o_j^{(t)}\| < 10^{-3} \cdot h_i$  then  

| Add  $o_j^{(t+1)}$  to set  $M$ ;  

| end  

| end  

| for each  $o_s^* \in M$  do  

|   for each  $o_s^* \in M, s \neq j$  do  

|     if  $\|o_j^* - o_s^*\| < \varepsilon$  then  

|       Calculate  $p_j = \frac{o_j^* + o_s^*}{2}$  and add  $p_j$  to set  

|        $P$ ;  

|     end  

|   end  

| end  

|  $P1 = unique(P)$   

|  $k = |P1|$ ;  

| if  $k > 0$  then  

|   Use elements in  $P1$  as initial cluster centers for  

k-means;  

|   Divide  $gb$  into  $k$  sub-GBs  

|    $\{sgb_1, sgb_2, \dots, sgb_k\}$ ;  

|   Add  $\{sgb_1, sgb_2, \dots, sgb_k\}$  to the end of  $Q$ ;  

| else  

|   Add  $gb$  to the end of  $GBs$  directly;  

| end  

| else  

|   Calculate: the center  $c$  and radius  $r$  of GB using  

Eq. (2) and Eq. (10)  

|    $GBs = GBs \cup gb$ ;  

| end  

| end  

38 Return: The GB set  $GBs$ .

```

The complete clustering procedure is formalized in Algorithm 2.

Algorithm 2: GB clustering algorithm

```

Input: GB set  $GBs = \{gb_1, gb_2, \dots, gb_m\}$ , radii  $r$ , centers  

 $c$ , initial local density peak set  $Modes$   

Output: Cluster labels  $CL$ 
1 Initialize:  $CL = \emptyset, GB_{nom} = \emptyset, GB_{out} = \emptyset, comp = \emptyset,$   

 $cluster = \emptyset$   

2 for  $gb_i \in GBs$  do  

3   if  $|gb_i| > 1$  then  

4      $GB_{nom} = GB_{nom} \cup \{gb_i\}$   

5   else  

6      $GB_{out} = GB_{out} \cup \{gb_i\}$   

7   end  

8 end  

9 for  $gb_i \in GB_{nom}$  do  

10   for  $gb_j \in GB_{nom}$  do  

11     Calculate distance  $d(gb_i, gb_j)$  by Eq. (21)  

12     Calculate normalized center distance  $s_{ij}$  using Eq.  

(22)  

13     Calculate similarity  $cd_{ij}$  by Eq. (23)  

14   end  

15 end  

16 for  $gb_i \in GB_{nom}$  do  

17   for  $mode \in Modes$  do  

18     Determine if  $gb_i$  is a core GB using Eq. (25)  

19   end  

20 end  

21 Let  $GBs_{core}$  be the set of core GBs  

22  $comp = BFS(cd_{ij}, GBs_{core})$   

23  $cluster = assignClusters(comp)$   

24 for  $gb_i \in GB_{out}$  do  

25    $cluster = cluster \cup \{gb_i\}$   

26 end  

27 for  $i = 1$  to  $|cluster|$  do  

28   if  $|cluster(i)| = 1$  then  

29     Let  $gb_s$  be the only GB in  $cluster(i)$   

30     Calculate distances from  $gb_s$  to every other GB in  

 $GBs$  using Eq. (21)  

31     Find the nearest GB  $gb_{nearest}$   

32     Assign the cluster label of  $gb_{nearest}$  to  $gb_s$   

33   end  

34 end  

35 for each data point  $x$  do  

36   if  $x$  is covered by some GB  $gb_k$  then  

37     Assign label of  $gb_k$  to  $x$   

38   else  

39     Calculate distances from  $x$  to all GBs by Eq. (26)  

40     Assign label of the nearest GB to  $x$   

41   end  

42 end  

43 Return: Cluster labels  $CL$  for all data points.

```

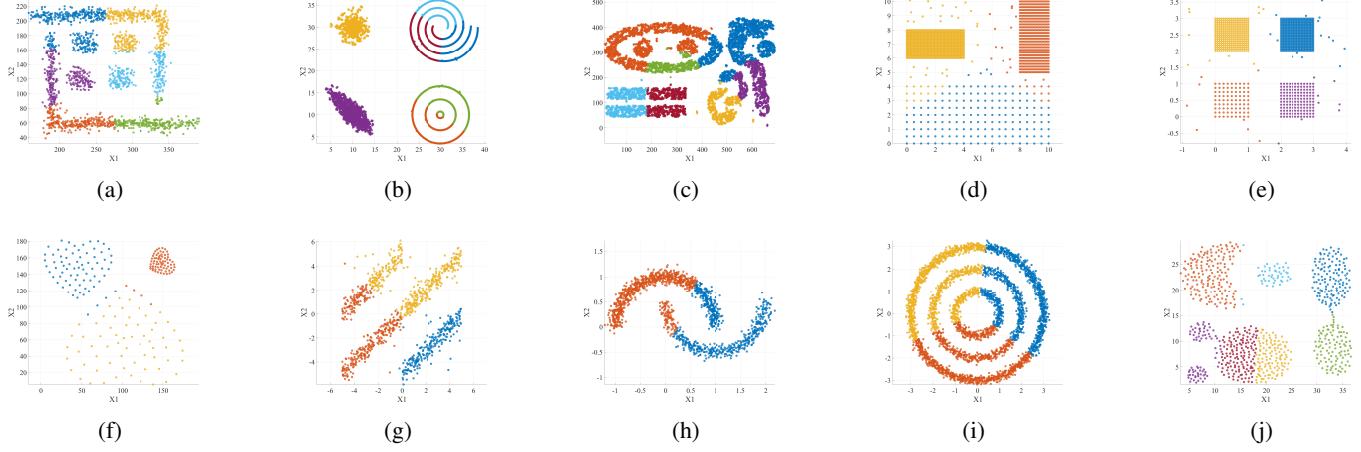


Fig. A.1: Clustering visualization of k-means on synthetic datasets

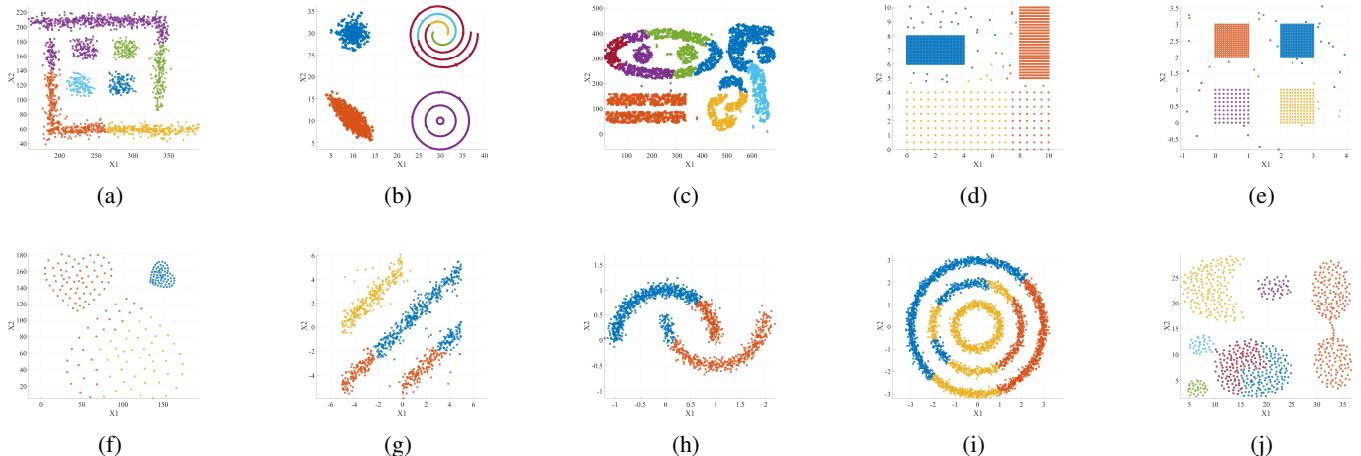


Fig. A.2: Clustering visualization of DP on synthetic datasets

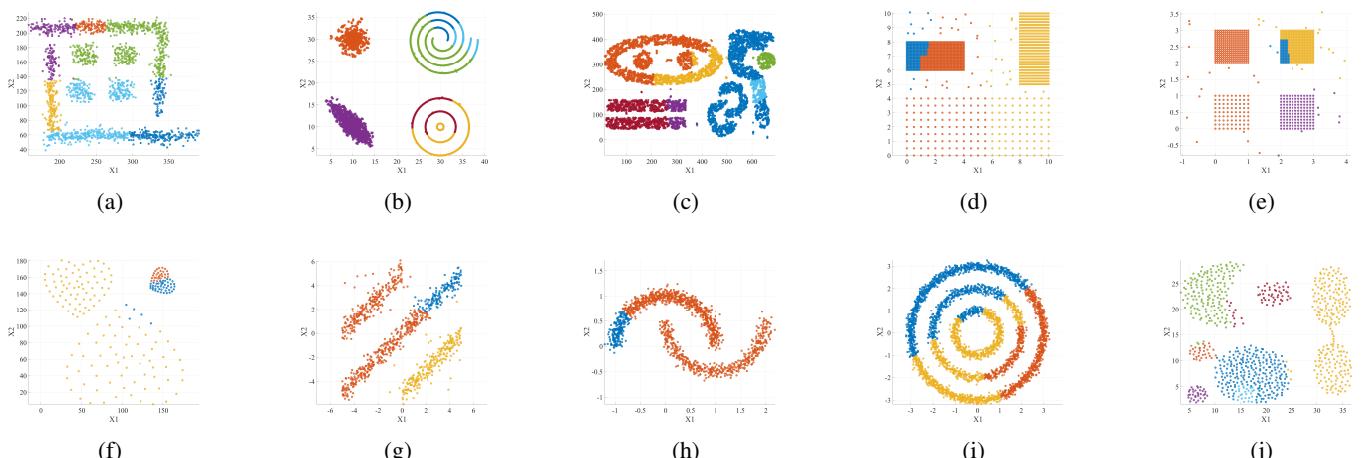


Fig. A.3: Clustering visualization of GB-DP on synthetic datasets

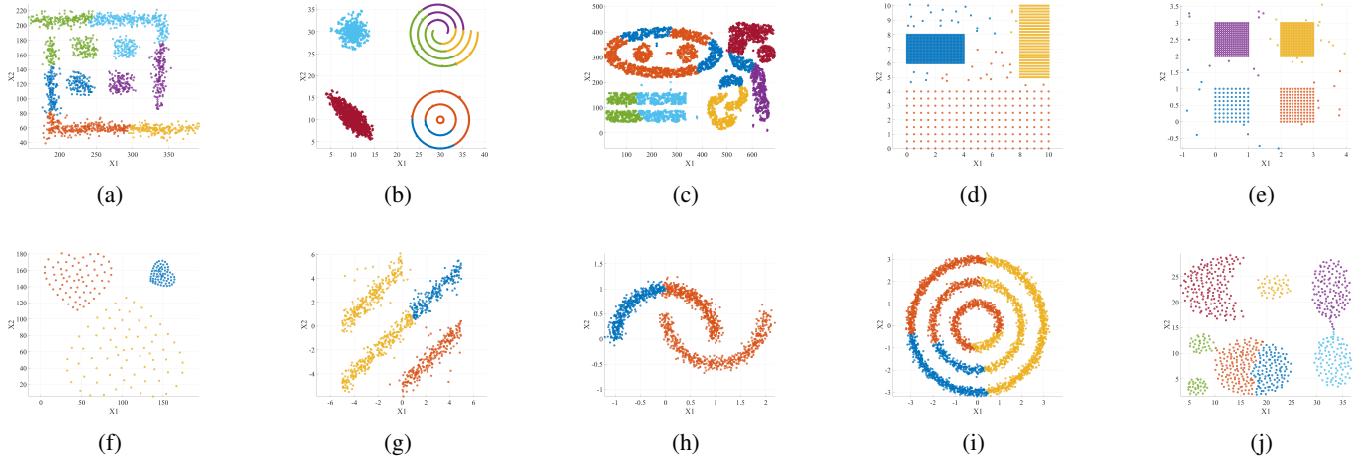


Fig. A.4: Clustering visualization of HC on synthetic datasets

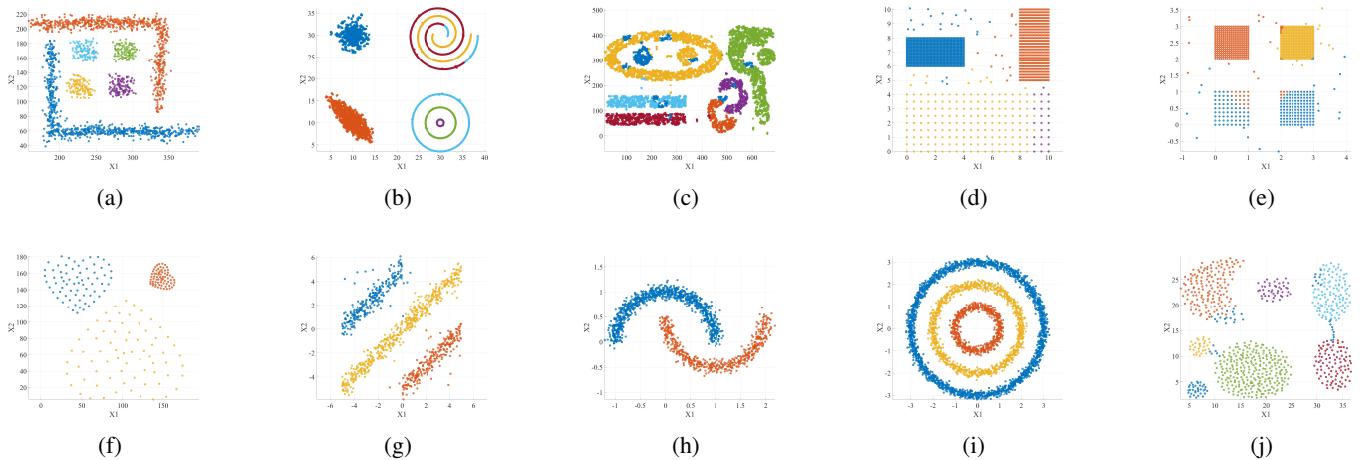


Fig. A.5: Clustering visualization of NaGB-DBSCAN on synthetic datasets

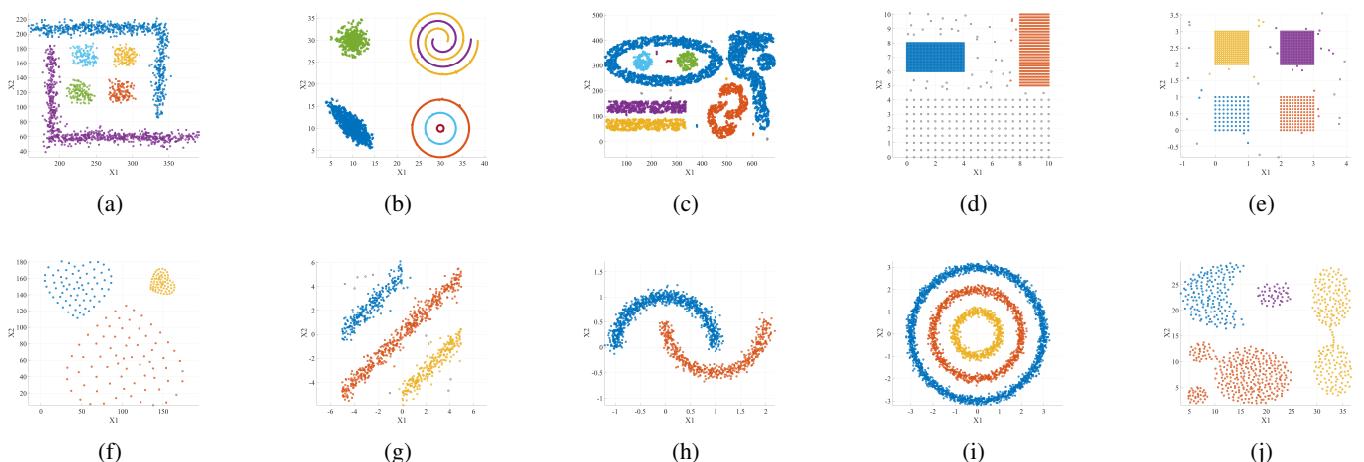


Fig. A.6: Clustering visualization of DBSCAN on synthetic datasets

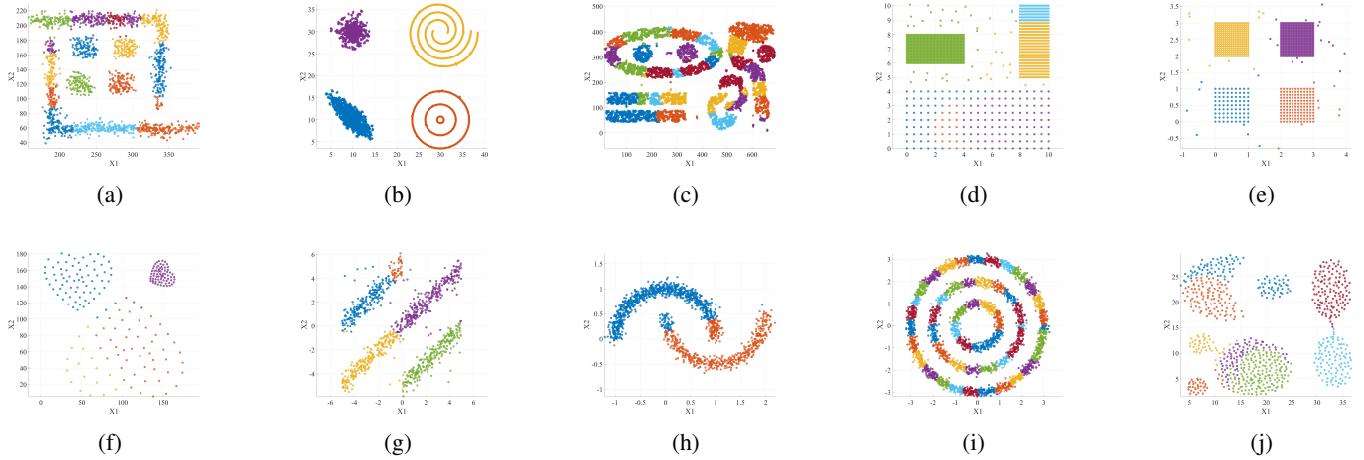


Fig. A.7: Clustering visualization of mean shift on synthetic datasets

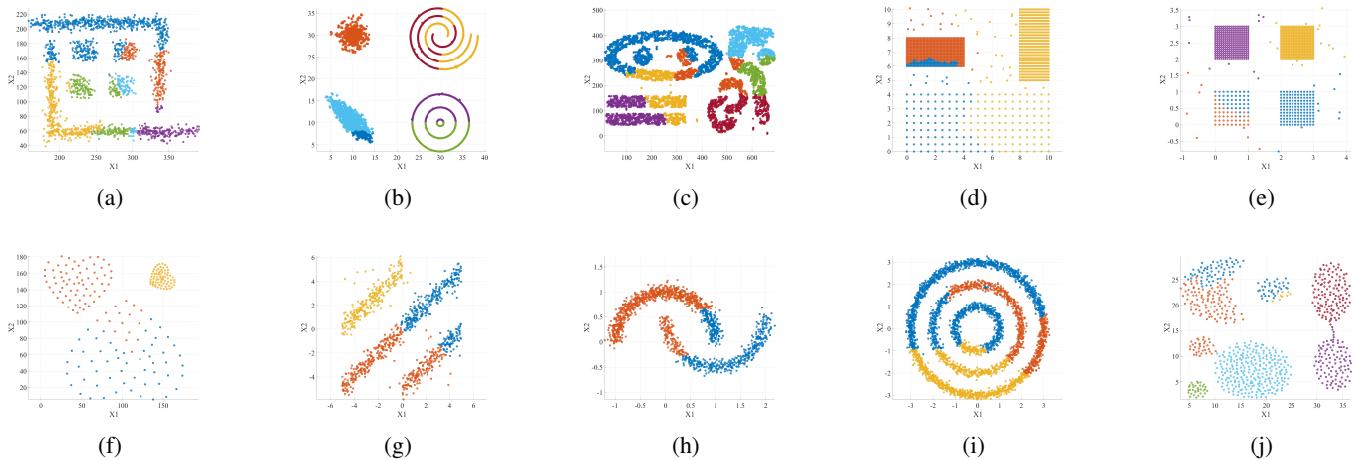


Fig. A.8: Clustering visualization of GBK-DPC on synthetic datasets

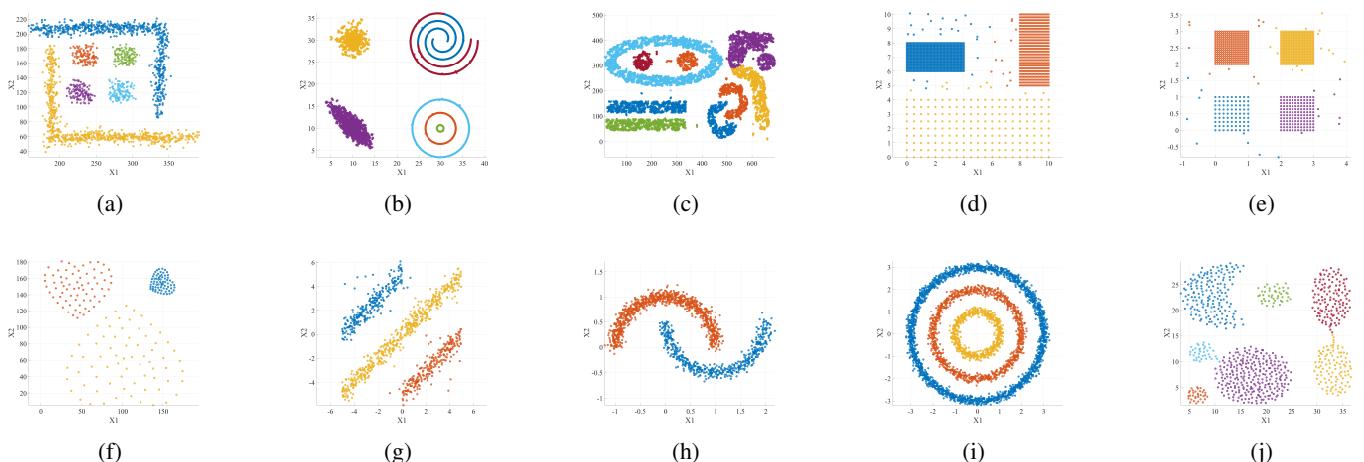


Fig. A.9: Clustering visualization of GBSC on synthetic datasets

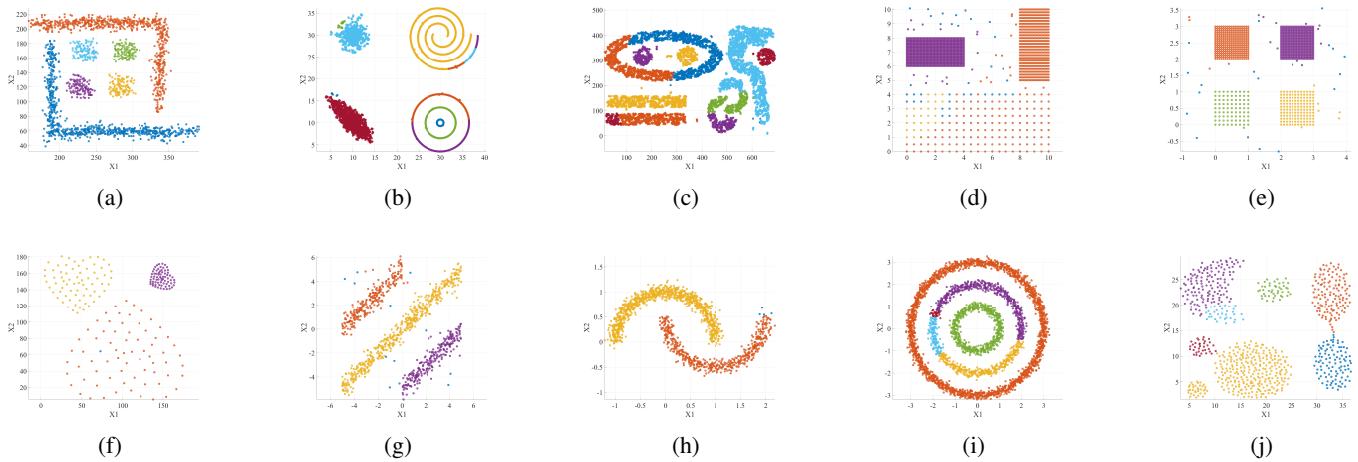


Fig. A.10: Clustering visualization of GBCT on synthetic datasets