

CURE

An Efficient Clustering Algorithm for Large Databases

presenter: Xiang Wan

Xiang Wan

The Introduction of Background

$$d_{mean}(C_i, C_j) = \|m_i - m_j\|$$

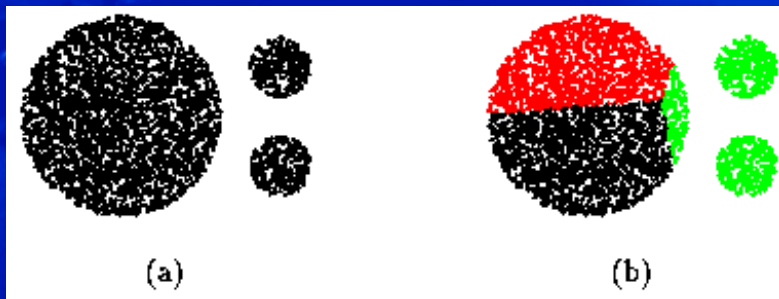
$$d_{ave}(C_i, C_j) = 1/(n_i n_j) \sum_{p \in C_i} \sum_{p' \in C_j} \|p - p'\|$$

$$d_{max}(C_i, C_j) = \max_{p \in C_i, p' \in C_j} \|p - p'\|$$

$$d_{min}(C_i, C_j) = \min_{p \in C_i, p' \in C_j} \|p - p'\|$$

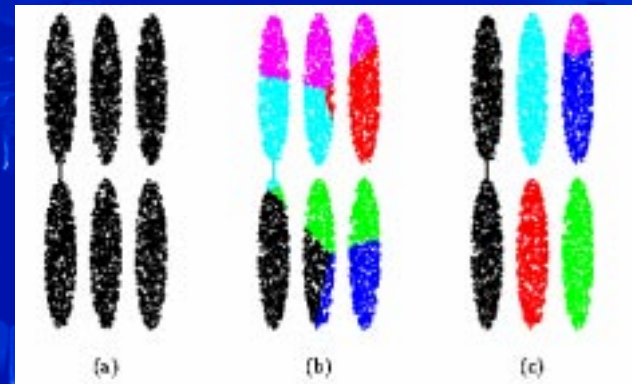
Xiang Wan

Example One



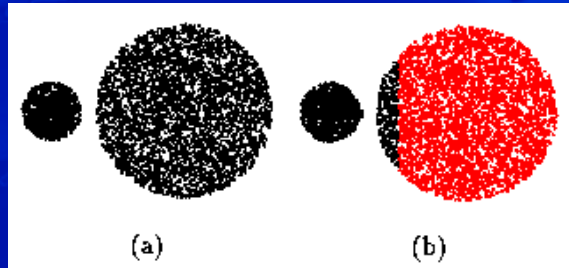
Xiang Wan

Example Two



Xiang Wan

Example Three



Xiang Wan

The Drawbacks of Traditional Clustering Algorithms

1. These algorithms can't identify the clusters having non-spherical shapes and wide variances in size.
2. These algorithms are sensitive to noises
3. These algorithms can't correctly label the data point when clusters don't have uniform sizes and shapes.

Xiang Wan

How to handle these problems in CURE

1. Choosing a constant number C of well scattered points in a cluster and using them to compute the distance between a pair of clusters.
2. Shrinking these points toward the centroid by a fraction α to mitigate the adverse effect of the noise.

Xiang Wan

Clustering Algorithms

```
procedure cluster(S,A)
begin
1. T := build_bstree(S)
2. Q := build_heap(S)
3. while size(Q) > 0 do {
4.   u := extract_min(Q)
5.   v := u.cluster
6.   delete(Q, u)
7.   w := merge(u, v)
8.   delete_max(T, u); delete_max(T, v); insert_max(T, w)
9.   u.cluster := w /* u is an arbitrary cluster in Q */
10.  for each x in Q do {
11.    if dist(x, u) < dist(x, u.cluster)
12.      u.cluster := x
13.    if u.cluster is either u or v {
14.      if dist(x, u.cluster) < dist(x, w)
15.        u.cluster := closest_cluster(T, x, dist(x, w))
16.    }
17.    u.cluster := w
18.    relocate(Q, x)
19.  }
20.  else if dist(x, u.cluster) > dist(x, w) {
21.    u.cluster := w
22.    relocate(Q, x)
23.  }
24. }
25. insert(Q, w)
26. }
end
```

Xiang Wan

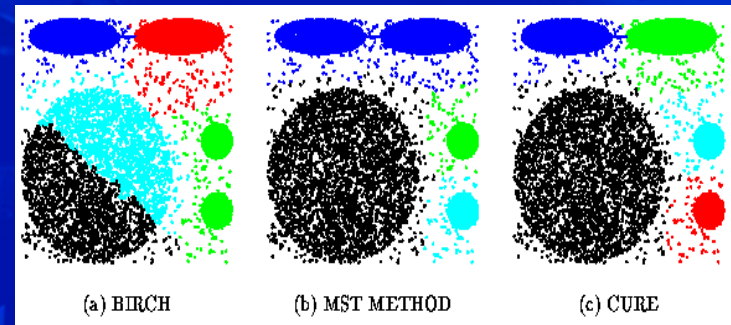
Merging Procedure

```

procedure merge(u, v)
begin
1. w := u ∪ v
2. w.mean :=  $\frac{u.mean + v.mean}{|u| + |v|}$ 
3. tmpSet := ∅
4. for i := 1 to c do {
5.   maxDist := 0
6.   foreach point p in cluster u do {
7.     if i = 1
8.       minDist := dist(p, w.mean)
9.     else
10.      minDist := min{dist(p, q) : q ∈ tmpSet}
11.     if (minDist ≥ maxDist){
12.       maxDist := minDist
13.       maxPoint := p
14.     }
15.   }
16.   tmpSet := tmpSet ∪ {maxPoint}
17. }
18. foreach point p in tmpSet do
19.   w.tmp := w.tmp ∪ {p + α(w.mean - p)}
20. return w
end
    
```

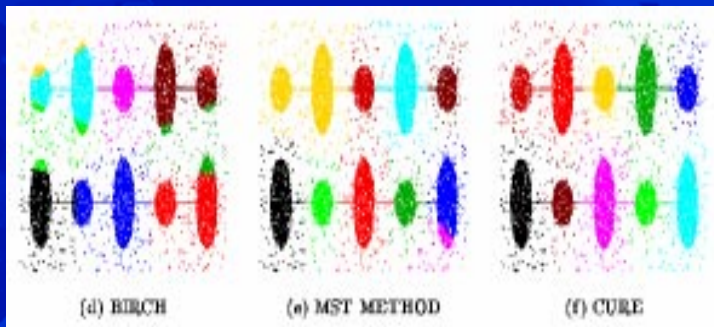
Xiang Wan

Example Four



Xiang Wan

Example Five



Xiang Wan