

, 1)

{yjsung, jpark}@cs.sungshin.ac.kr

가

가

MST, CLARANS, CURE PROCLUS
CLARANS PROCLUS
PROCLUS

가

2 3

가

confusion matrix

1.

가

가

[8, 13, 14].

가

MST, CLARANS[2], CURE[1]

PROCLUS[3]

가 PROCLUS CLARANS

3

2

가

2. , 3 , 4 , 5 .

2.

2.1

2.1.1 MST (Minimum Spanning Tree)

Hierarchical 가 가 . Ouliter .
 k 가 chaining effect가 .
 outlier가 .

2.1.2 CLARANS (Clustering Large Application based on Randomized Search)

CLARANS PAM CLARA CLARA가
 CLARANS
 k-medoid
 k-medoid medoid
 (neighbor) , maxneighbor
 CLARANS 가
 k-medoid
 (local optimal) CLARANS
 numlocal [2].

2.1.3 CURE (Clustering Using Representatives)

k 가
 가 hierarchical
 c
 가 c
 shrinking factor
 c outlier chaining effect
 MST
 centroid-based approach 가

2.2

(sparsity) 가
 가 가 가
 . [1] 3 x-y
 3
 가 2
 가 [2] 3
 가

2, k [0.0, 1.0)

outlier 5%

3.1.1

[3]

3.2

MST, CLARANS, CURE

PROCLUS

PROCLUS

medoid greedy

A • k A > B B • k medoid greedy

outlier가 greedy

A • k B • k medoid k

medoid medoid M_{current}가

M_{current} locality medoid

PROCLUS 가 k medoid

$M = \{m_1, \dots, m_k\}$ m_i $\delta_i = \min_{j=1, \dots, k} d(m_i, m_j)$ δ_i [3] k가

3 medoid, $M = m_1, m_2, m_3$ δ_i 2 medoid

locality L_i L_i m_i δ_i

$(L_1, \dots, L_k$ disjoint $)$ L_i m_i

δ_i medoid m_i 가 가 j

medoid k • l Medoid

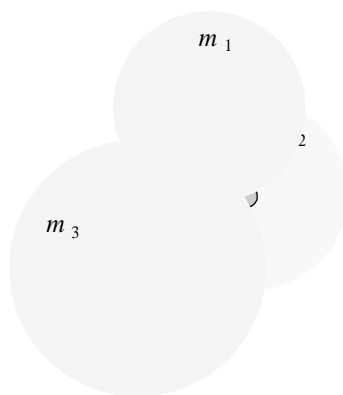
medoid 가 가 medoid

Medoid medoid M_{current}

가 M_{current} 가 M_{best} M_{current}가

M_{current} M_{best} termination_criterion

M_{best} locality



[3] Locality

4.

가

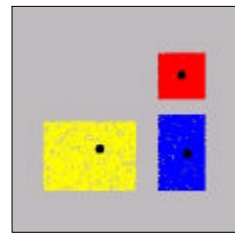
confusion matrix

4.1 Outlier

outlier
10000, outlier
5%
MST가 outlier가, CLARANS CURE
outlier
[5] MST, CLARANS, CURE
outlier가 [4]

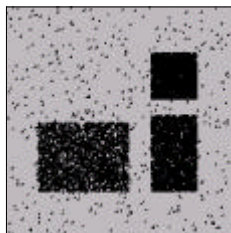


(a)

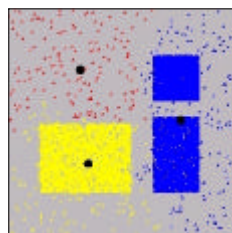


(b) MST, CLARANS, CURE

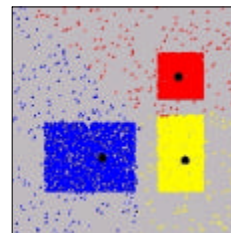
[4] Outlier가



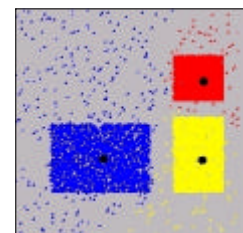
(a)



(b) MST



(c) CLARANS



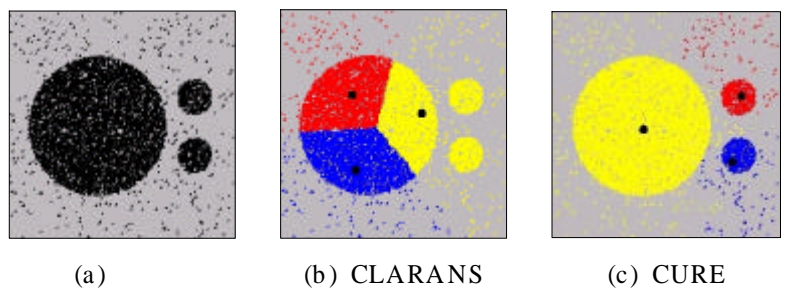
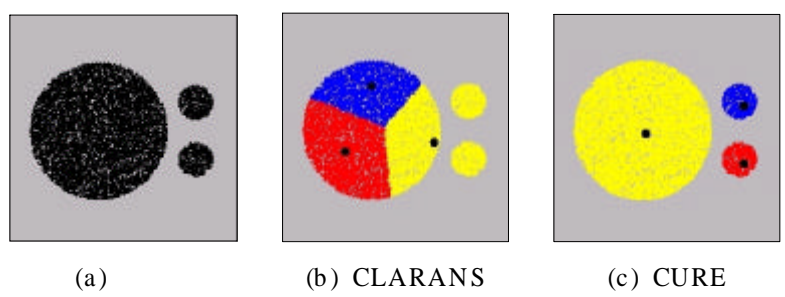
(4) CURE

[5] 5% outlier 가

4.2

[6] [7]
(
4),
CURE
outlier
[6] 5%
CLARANS
outlier가
CLARANS
10000

centroid-based
 CURE
 $c = 10$, shrinking factor = 0.5
 5% outlier



4.3
 CLARANS PROCLUS
 20 10000, outlier
 4
 [2] CLARANS
 가

[4] PROCLUS
 가
 [5] confusion matrix
 가 가 outlier outlier

[8] [9]
 3
 confusion matrix
 [8] [9] (b)PROCLUS 4 가
 , [8] [9] (a)CLARANS
 가 가

Input	Dimensions	Points
A	2, 4, 11, 14	2584
B	2, 3, 7, 14	2812
C	2, 12, 14, 17	2204
D	2, 3, 13, 14	1900
Outliers		500

[1] $n = 10000$, $k = 4$, $\epsilon = 4$

Found	Dimensions	Points
1	Full dimensions	2701
2	Full dimensions	1736
3	Full dimensions	1965
4	Full dimensions	3598
Outliers		

[2] CLARANS

Input	A	B	C	D	Out.
Output					
1	834	1062	21	649	135
2	283	276	156	953	68
3	196	5	1602	25	137
4	1271	1469	425	273	160
Outliers					182

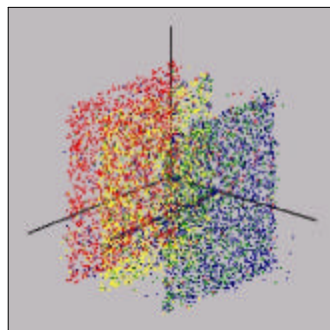
[3] CLARANS: Confusion Matrix

Found	Dimensions	Points
1	2, 4, 11, 14	2659
2	2, 3, 7, 14	2859
3	2, 12, 14, 17	2360
4	2, 3, 13, 14	1940
Outliers		182

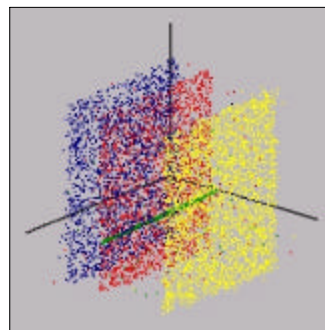
[4] PROCLUS

Input Output	A	B	C	D	Out.
1	0	0	0	1900	135
2	0	2812	0	0	47
3	2584	0	0	0	75
4	0	0	2204	0	156
Outliers					182

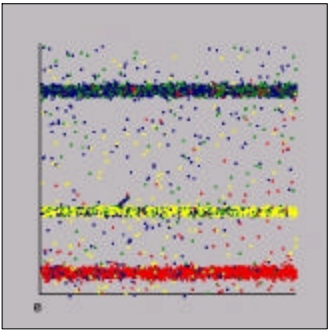
[5] PROCLUS: Confusion Matrix



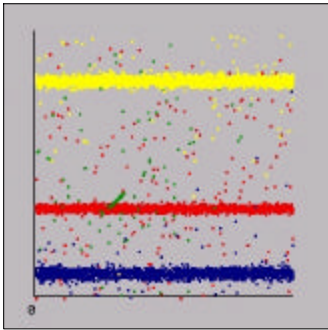
(a) CLARANS



(b) PROCLUS



(a) CLARANS



(b) PROCLUS

[9] 20

2

5.

가

가

[9, 15].

가

가

k

5.1 2

k

. Centroid-based

center

medoid

[10] 2

Display Data

5.2

3

d

2

3

d

,

2

3

Display

k

3

3

[11] 3

[12]

20

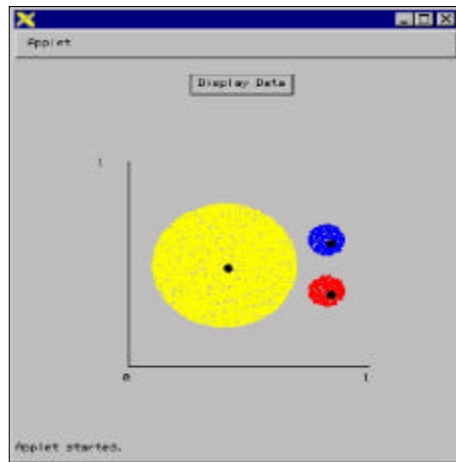
. [13]

[12]

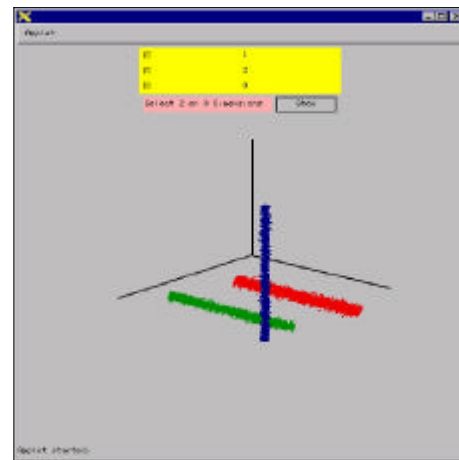
3

2

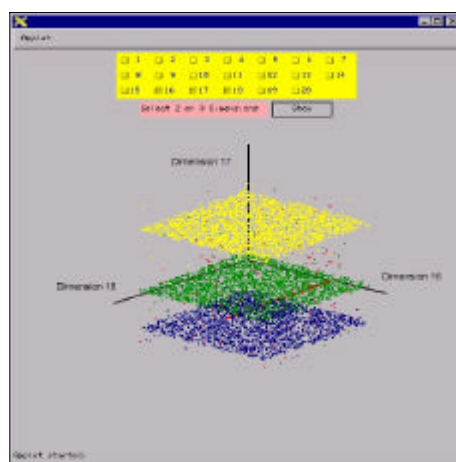
2



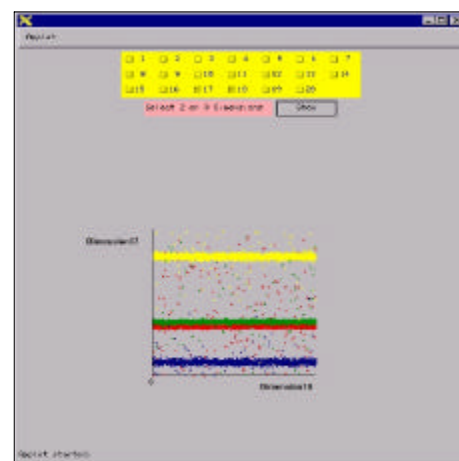
[10] 2



[11] 3



[12] 20
3



[13] 20
2

6.

가 가

가

(MST, CLARANS, CURE) CURE가 PROCLUS 가

가

가

boolean categorical

- [1] S. Guha, R. Rastogi and K. Shim, CURE: "A Efficient Clustering Algorithm for Large Databases", In *Proceedings of A CM SIGMOD*, pages 73-84, 1998.
- [2] R. Ng and J. Han, "Efficient and Effective Clustering Method for Spatial Data Mining", In *Proceedings of the 20th VLDB Conference*, pages 144- 155, 1994.
- [3] C. C. Aggarwal, C. Procopiuc, J. L. Wolf, P. S. Yu, and J. S. Park, " Fast Algorithms for Projected Clustering," In *Proceedings of the A CM SIGMOD International Conference on Management of Data*, pp. 61-72, Philadelphia, PA, June 1-3, 1999.
- [4] R. Agrawal, J. Gehrke, D. Cunoploes, P. Raghavan, "Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications", In *Proceedings of the A CM SIGMOD International Conference on Management of Data*, 1998.
- [5] S. Guha, R. Rastogi, and K. Shim, "ROCK: A Roust Clustering Algorithm for Categorical Attributes", the 15th International Conference on IEEE Data Engineering, 1999.
- [6] M. Ester, H.-P. Kreigel, J. Sander, "A Density Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise", In *Proceedings of the 2nd International Conference on Knowledge Discovery in Databases and Data Mining*, 1995.
- [7] T. Zhang, R. Ramakrishnan, M. Livny, "BIRCH: An Efficient Data Clustering Method for Very Large Databases", In *Proceedings of the A CM SIGMOD International Conference on Management of Data*, 1996.
- [8] R. Agrawal, M. Mehta, J. Shafer, and R. Srikant, "The Quest Data Mining System", In *Proceedings of the 2nd International Conference on Knowledge Discovery in Databases and Data Mining*, 1996.

- [9] T. Fukuta, Y. Morimoto, and S. Morishita, "Data Mining Using Two-Dimensional Optimized Association Rules: Scheme, Algorithms, and Visualization", In *Proceedings of A CM SIGMOD*, pages 13-23, 1996.
- [10] R. Agrawal, T. Imielinski and A. Swami, "Mining Association Rules between Sets of Items in Large Databases", In *Proceedings of A CM SIGMOD*, pages 207-216, 1993.
- [11] R. Agrawal and R. Srikant, "Mining Sequential Patterns", In *Proceedings of the 11th International Conference on Data Engineering*, pages 3-14, 1995.
- [12] M.-S. Chen, J. S. Park, and P. S. Yu, "Data Mining for Path Traversal Patterns in a Web Environment", In *Proceedings of the 16th International Conference on Distributed Computing Systems*, pages 385-392, 1996.
- [13] M. S. Chen, J. Han, and P. S. Yu, "Data Mining: An Overview from Database Perspective", *IEEE Transactions on Knowledge and Data Engineering*, Vol. 8, No. 6, pp. 866-883, Dec. 1996.
- [14] R. Agrawal, T. Imielinski and A. Swami, "Database Mining: A Performance Perspective", *IEEE Transformation on Knowledge and Data Engineering*, pages 914-925, 1993.
- [15] S. M. Weiss and N. Indurkha, "Predictive Data Mining", Morgan Kaufmann Publishers, Inc, 1998.
- [16] , “ ”,
 , 1999.
- [17] , “ ”,
 , 1997.