

가

(Design and development of the clustering algorithm  
considering weight in spatial data mining)

(Ho-Sook Kim)

(Hyun-Sook Lim)

(Hwan-Seung Yong)

가

DBSCAN-W

DBSCAN-W

DBSCAN

DBSCAN

, DBSCAN-W

DBSCAN-W

가

DBSCAN-W

**Abstract**

Spatial data mining is a process to discover interesting relationships and characteristics those exist implicitly in a spatial database. Many spatial clustering algorithms have been developed. But, there are few approaches that focus

simultaneously on clustering spatial data and assigning weight to non-spatial attributes of objects. In this paper, we propose a new spatial clustering algorithm, called DBSCAN-W, which is an extension of the existing density-based clustering algorithm DBSCAN. DBSCAN algorithm considers only the location of objects for clustering objects, whereas DBSCAN-W considers not only the location of each object but also its non-spatial attributes relevant to a given application. In DBSCAN-W, each datum has a region represented as a circle of various radius, where the radius means the degree of the importance of the object in the application. We showed that DBSCAN-W is effective in generating clusters reflecting the user's requirements through experiments.

: , , 가

1.

[1].

(partitioning) [2] (hierarchical) [3,4], (density-based)  
 [5,6], (grid-based) [7] (model-based) [8].

가

가

[8][9].

가

가

1 : 가 automatic teller machine (ATM)

가

가

(noise)

가

1

가

DBSCAN - W

DBSCAN - W

DBSCAN

DBSCAN

, DBSCAN - W

DBSCAN

가

, DBSCAN - W

가

가

가

DBSCAN - W가

. 2

DBSCAN

, 3

DBSCAN - W

가

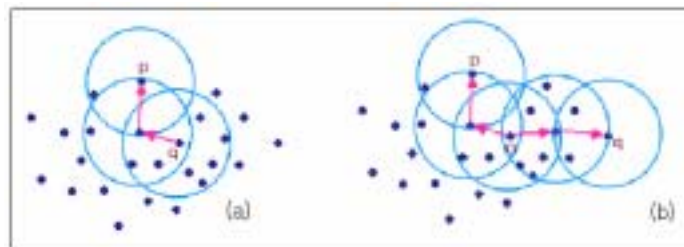
, 4

2. DBSCAN :

[5]

DBSCAN(Density Based Spatial Clustering of Applications with Noise)[5]

- DBSCAN 가
- $p$  Eps-neighborhood  $p$  Eps (neighborhood)
- MinPts(minimum number of points) Eps-neighborhood 가 core object
- $D$   $p$ 가  $q$  directly density-reachable  $p$ 가  $q$  Eps-neighborhood  $q$ 가 core object
- $p$ 가  $q$  density-reachable  $p$   $q$  directly density-reachable (chain)
- $p$ 가  $q$  density-connected  $p$   $q$  density-reachable
- $p$ 가  $q$  (density-connected)



( 1) Density-reachability density connectivity  
 (a)  $p$   $q$  density-reachable

(b)  $\sigma, p, q$  density-connected

DBSCAN 1 Eps  
Eps  
(overlap)  $\log n$  R\*-tree  
DBSCAN 가  $n$   
 $O(n * \log n)$

### 3. 가

3 가  
DBSCAN DBSCAN-W ( a DBSCAN algorithm using region expressed as Weight)  
. DBSCAN  
. , DBSCAN-W

DBSCAN-W 가  
1. 가  
가  
2.  $p$  Eps-neighborhood  $p$  Eps  
3. - (maximal set of density-connected regions)

DBSCAN-W

3

가

A

$$F(A_i) = r_i$$

( $r_i$ )

2

( $r_i$ )

DBSCAN-W가

가

2 가

F

가

N

[0, N]

가

[0,60]

가

가

10

가 1, 10

100

가 2

가

가

가

가 ' '

가 10, ' '

가 1

DBSCAN

Core-point

p

seed

,

seed

Eps

MinPts

density-reachable

DBSCAN-W

core-point

p

Eps

MinPts

density-reachable

Eps-neighbor

2

2

DBSCAN

DBSCAN-W

Eps-neighborhood

2(a)

2(b) DBSCAN neighborhood  $p$  Eps-neighborhood  $p$  Eps-neighborhood

neighborhood  $p$  Eps 5 2(c)

DBSCAN-W Eps-neighborhood

(x,y) r

가 가

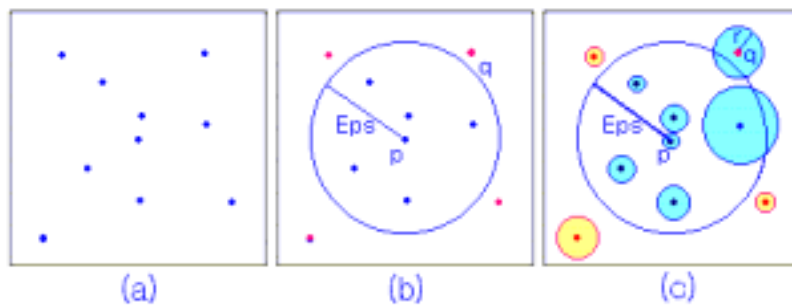
가

2 (c) q

p Eps p Eps q

p Eps-neighborhood 가 가

가



( 2) DBSCAN DBSCAN-W  $p$  Eps-neighborhood

(a)

(b) DBSCAN  $p$  Eps-neighborhood

(c) DBSCAN-W  $p$  Eps-neighborhood

4.

3 DBSCAN-W 가 1

ATM

DBSCAN

(synthetic)

Sun

Solaris 2.6

-

DBMS Informix Universal Server

Informix Spatial Datblade Module 2

9

가

R-

Tree [10,11].

Create table Customer

```

( Customer_id      smallint not null, /* ID */
  Location         sp2Pnt, /* (x, y) */
  Region          sp2Circ, /* 가 (x, y, r) */
  Avg_trade_frequency smallint, /* */
  Avg_trade_amount smallint, /* */
  Weight          smallint, /* 가 */
  Cluster_id      smallint /* */
);

```

1.

7

500

Customer\_id . Location (x, y)

1 ~ 1000

. Region location ,

Weight type . Avg\_trade\_frequency Avg\_trade\_amount

DBSCAN-W

가

. Weight

가 가

F

Region

Cluster\_id

Eps

MinPts

가



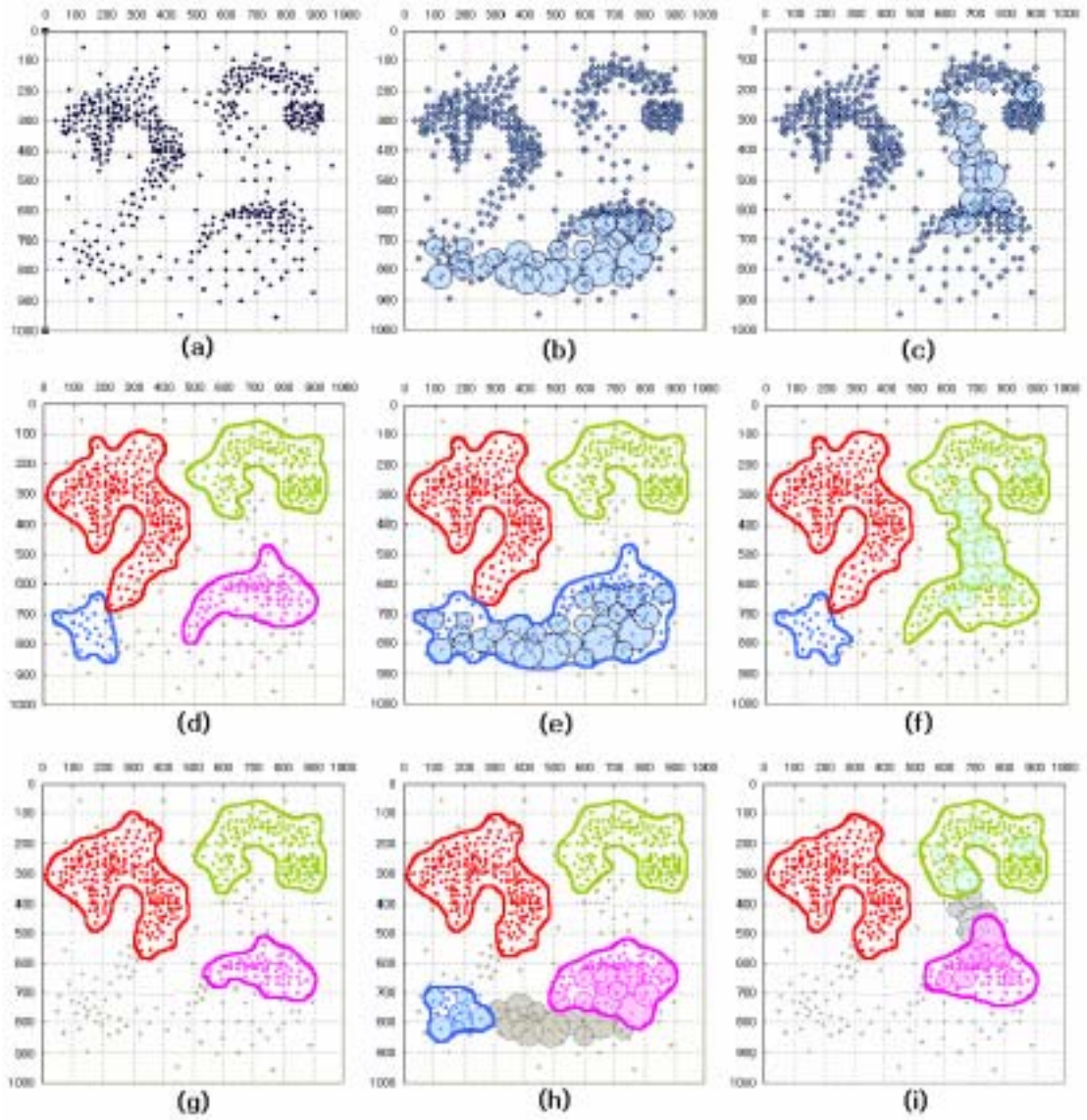
3가

F 가

[0,60]

DBSCAN DBSCAN-W

3



( 3 )

(a) Location

(b) Avg\_trade\_frequency

(c) Avg\_trade\_amount

(d) Eps =50, MinPts =5

가

가

DBSCAN

Region

Region

- (e) Eps =50, MinPts =5, 가 = Avg\_trade\_frequency DBSCAN-W
- (f) Eps =50, MinPts =5, 가 = Avg\_trade\_amount DBSCAN-W
- (g) Eps =50, MinPts =10 DBSCAN
- (h) Eps =50, MinPts =10, 가 = Avg\_trade\_frequency DBSCAN-W
- (i) Eps =50, MinPts =10, 가 = Avg\_trade\_amount DBSCAN-W

3 (a) 2 500 가 3 (b)

(c) 가 Avg\_trade\_frequency Avg\_trade\_amount

가 region 가

3 (b)

3 (c) 3 (d) Eps =50, MinPts =5

DBSCAN 4 가 ,

noise 3 (e) Eps =50, MinPts =5

가 DBSCAN-W

가 3 (d)

DBSCAN

3 (f) Eps =50,

MinPts =5 가 DBSCAN-W

가

3 (g)~(i)

DBSCAN DBSCAN-W

MinPts 가 3 (d)~(f)

가 가 가

가

가 DBSCAN-W DBSCAN

5.

가

DBSCAN-W

DBSCAN

가

가

[1] Michael J. A Berry, and Gordon Linoff, Data Mining Techniques : For Marketing, Sales, and Customer Support, John Wiley & Sons, Inc., 1997.

[2] Raymond T. Ng, and Jiawei Han, "Efficient and Effective Clustering Method for Spatial Data Mining," In Proc. of the VLDB Conference, Santiago, Chile, pp. 144-155, September 1994.

[3] Tian Zhang, Raghu Ramakrishnan, and Miron Livny, "BIRCH : An Efficient Data Clustering Method for Very Large Databases," In Proc. of the ACM SIGMOD Conference on Management of Data, Montreal, Canada, pp. 103-114, June 1996.

[4] Sudipto Guha, Rajeev Rastogi, and Kyuseok Shim, "CURE : An Efficient Clustering

- Algorithm for Large Databases," In Proc. of the ACM SIGMOD Conference on Management of Data, Seattle, Washinton, USA, pp. 73-84, May 1998.
- [5] Martin Ester, Hans-Peter Kriegel, Jorg Sander, and Xiaowei Xu, "A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise," In Proc. of ACM SIGMOD 3rd International Conference on Knowledge Discovery and Data Mining, pp. 226-231, AAAI Press, 1996.
- [6] Mihael Ankerst, Markus M. Breuning, Hans-Peter Kriegel, and Jorg Sander, "OPTICS: Ordering Points To Identify the Clustering Structure," In proc. of ACM SIGMOD International Conference on Management of Data, Philadelphia, Pennsylvania, USA, pp. 49-60, June 1999.
- [7] W.Wang, J.Yang, and R.Muntz, "STING :A statistical information grid approach to spatial data mining", In Proc. 1997 Int. conf. Very Large Data Bases(VLDB'97), Athens, Greece, pp.186-195, August 1997.
- [8] Jiawei Han, and Micheline Kamber, Data Mining: Concepts and Techniques, Morgan Kaufmann publishers, 2001.
- [9] Douglass R. Cutting, David R. Karger, Jan O. Pedersen, and John W. Tukey, "Scatter/Gather: A Cluster-based Approach to Browsing Large Document Collections," In Proc. of the 15th Annual International ACM SIGIR Conference, pp. 318-329, June 1992.
- [10] Informix, Informix Universal Server Guide to SQL: Tutorial Version 9.1, Informix Press, 1997.
- [11] Informix, Informix Spatial Datablade Module: User's Guide, Informix Press, 1997.

(Ho-Sook Kim)

1993 :  
1993 ~1997 : SDS  
1999 :  
1999 ~ :  
2001 ~ :  
< > ,

(Hyun-Sook Lim)

1999 :  
2002 :  
2002~ : LG CMS  
< > , CRM

(Hwan-Seung Yong)

1983 :  
1985 :  
1985 ~ 1989 :  
1994 :  
1995 ~ :  
< > , XML