

2001

가

2002

가

論文 碩士學位 論文 提出

2002 年 1月

梨花女子大學校 科學技術大學院

學科 林 賢 淑

碩 士 學 位 論 文 認 准

指導教授

審查委員

_____	_____
_____	_____
_____	_____

梨 花 女 子 大 學 校 科 學 技 術 大 學 院

	-----	v
I.	-----	1
1.1	-----	1
1.2	-----	2
II.	-----	4
2.1	-----	4
2.1.1	-----	4
2.1.2	-----	5
2.1.3	- DBSCAN -----	6
2.2 가 (Visibility Graph)	-----	9
III. 가	-----	13
3.1	가 -----	13
3.1.1	-----	14
3.1.2 가	-----	14
3.2	-----	18
3.2.1	-----	19
3.2.2	-----	19

IV.	가	-----	23
4.1		-----	23
4.2		-----	25
4.3	가	-----	27
4.3.1		-----	27
4.3.2		-----	29
4.4	가	-----	29
V.		-----	35
		-----	36
		-----	39

[4.1]	-----	24
[4.2]	-----	25
[4.3]	-----	26
[4.4]	-----	29

論文概要

DBSCAN

가

가

DBSCAN 가

DBSCAN-W COD-DBSCAN DBSCAN-W

가

, COD-DBSCAN

가

I.

1.1

(clustering) 가 , 가 [1].

(data mining) , , 가 (, I/O) 가 , 가 .

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

,

[5].

가

가 ,

1.2

가

가

DBSCAN-W

direct

Euclidian distance

COD-DBSCAN

가

DBSCAN

가

. II

DBSCAN

가

(Visibility Graph)

가

. III

DBSCAN-W

COD-DBSCAN

, IV

가

.

V

.

II.

(Spatial Data)

DBSCAN

(robot motion planning)

가

2.1

2.1.1

(spatial data)

(discrete)

(continuous)

(point)

(distance

attribute)

(non-spatial data)

(region)

[9].

2.1.2

가

가

[10].

가

[10].

1) (input parameter)

[11].

2) (ball-shaped),

(arbitrary shape)

(size)

[5,11].

3)

(noise)

[5].

4)

[5,11].

2.1.3

- DBSCAN

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

가 , (noise)
 (arbitrary shape) (size)
 (density)

[1] (Eps-neighborhood of a point) p Eps-neighborhood p Eps
 (neighborhood)

$$N_{Eps}(p) = \{ q \in D \mid \text{dist}(p, q) \leq Eps \}$$

[2] (directly density-reachable) Eps MinPts가 p

q directly density-reachable

1) $p \in N_{Eps}(q)$

2) $|N_{Eps}(q)| \geq \text{MinPts}$ (core point)

[3] (density-reachable) p q directly density-reachable

(chain) p q density-reachable

[4] (density-connected) p q density-reachable 오가

p q density-connected .

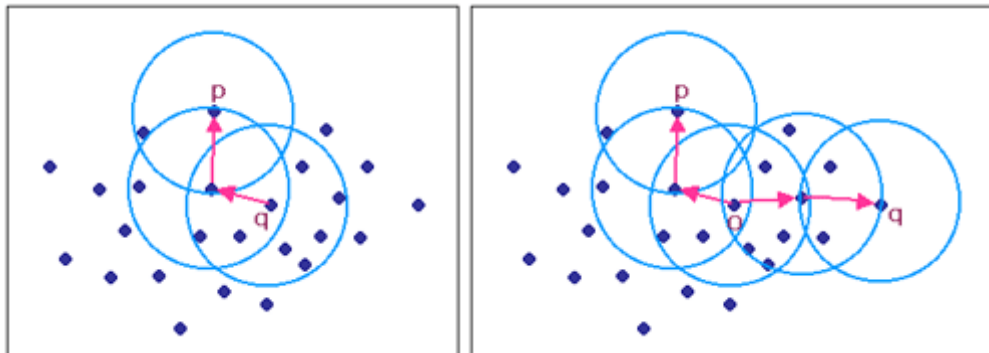
[5] (cluster) D , C

D non-empty subset .

1) p, q $p \in C$ $q \notin p$ density-reachable $q \in C$

2) $p, q \in C$ p q density-connected .

[6] (noise)



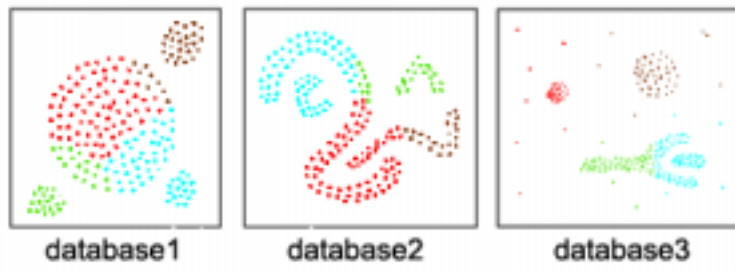
(a) density-reachable

(b) density-connected

[2.1] density-reachable density-connected

2.1 (a) $p \notin q$ density-reachable , (b) p

q 가 o density-connected . ,
 - (density-connected)
 Eps MinPts . R*-
 Tree , DBSCAN 가
 n $O(n * \log n)$



(a) CLARANS (k = 4)



(b) DBSCAN

[2.2] CLARANS DBSCAN

2.2 가 4 (ball-shaped)
 1, (non-convex shape) 2,

(arbitrary shape)

3

,

CLARANS

DBSCAN

2.2 가 (Visibility Graph)

가

(Visibility Graph)

가

(robot motion planning)

,

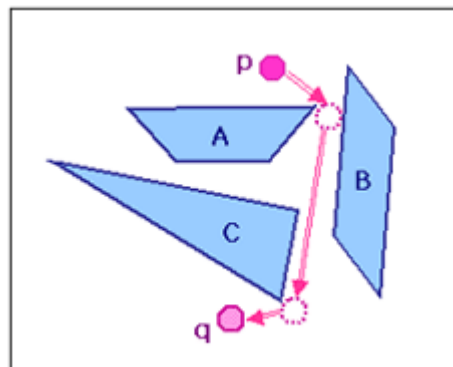
2.3

p

q

A, B, C

[12].



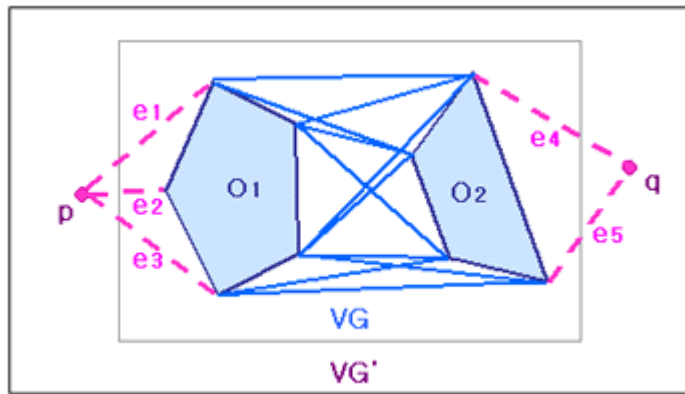
[2.3]

node V , (visible)

edge E 가 VG

$$VG = (V, E)$$

2.4 가 가 , 가
 가 . O_1, O_2 가 $VG = (V,$
 E) , p q 가 .
 $V' = V \cup \{p, q\}$ $E' = E \cup \{e_1, e_2, e_3, e_4, e_5\}$,
 $VG' = (V', E')$



[2.4] 가

가 naïve
 가 $O(n^3)$.
 edge (cyclic order) 가
 $O(n^2 * \log n)$,
 2.5 [12,13]. 2.6 (a)
 , edge (b) .

VisibilityGraph()

for all vertices w_i where $i = 1$ to n

if **VISIBLE**(w_i) **then Add** w_i to the list of visibility edges

Insert into T the obstacle edges incident to w_i that lie on the

clockwise side of the half-line from p to w_i

Delete from T the obstacle edges incident to w_i on the

counterclockwise side of the half-line from p to w_i

VISIBLE(w_i)

if pw_i intersects the interior of the obstacle of which w_i is a vertex

then return false

else if $i = 1$ or w_{i-1} is not on the segment pw_i

then Search in T for the edge e in the leftmost left

(the edge p intersects first)

if e exists and pw_i intersects e then return false

else return true

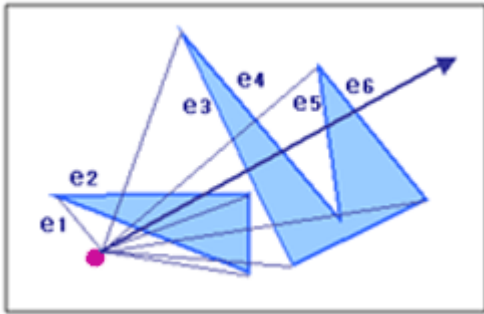
else if w_{i-1} is not visible **then return false**

else Search in T for an edge e that intersects $w_{i-1}w$

if e exists then return false

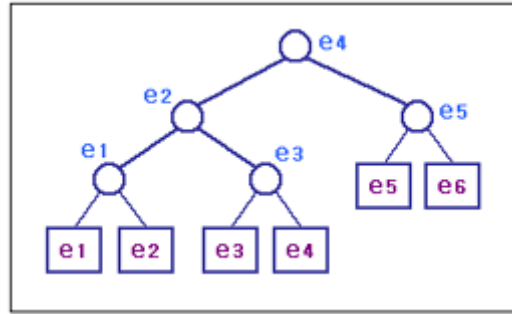
else return true

[2.5] 가



(a)

[2.6]



(b)

edge

edge

III. 가

가

가

3.1

가

가

가

가

가

가

DBSCAN-W

3.1.1

[3.1]

가 가
가

[3.2]

(, ,
가) 가

3.1.2 가

가

4가

, 가 가 가 (Eps)

, 가

가

, 가 가 가 (MinPts)

, 가

가 DBSCAN (global parameter)

가

, DBSCAN 가

(location) .

, 가 가

, .

3.1 가

가 . 3.1 (a) 가

가 가 1

가 가 2 . (b)

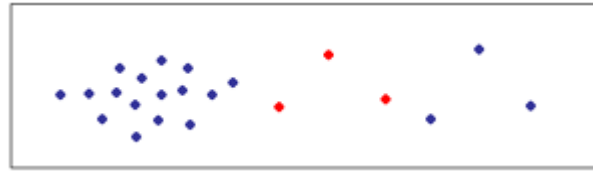
가

(c) . (d) , 가

, 가 .

(e) . (overlap) 가

.



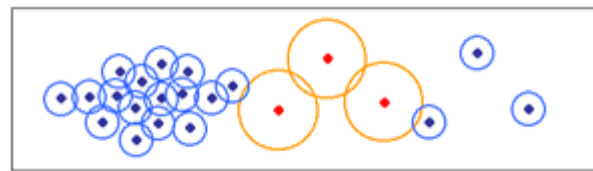
(a) 가 가



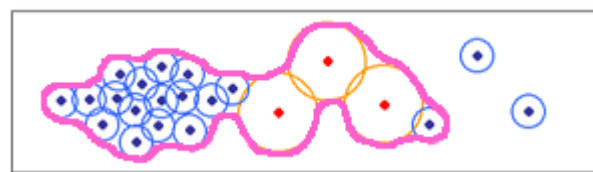
(b) 가



(c) 가



(d) 가



(e) 가

[3.1] (overlap)
가

가

DBSCAN-W(DBSCAN algorithm

considering Weight)

가

DBSCAN

가

가

[1]

(x, y)

, 가

[2]

p Eps-neighborhood p

Eps

[3]

(maximal set of density-connected

circles)

3.2 DBSCAN-W

DBSCAN

(a)

, (b)

가

DBSCAN

6

p

Eps

6

p Eps-neighborhood . (c)

가

DBSCAN-W

, (b)

, p Eps

가 가 q가 가

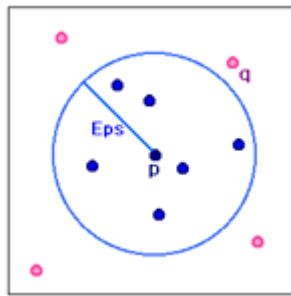
p

Eps-neighborhood . , 가

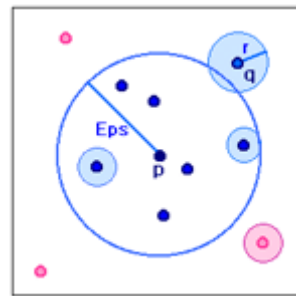
가 가



(a)



(b) DBSCAN



(c) DBSCAN-W

[3.2] DBSCAN DBSCAN-W

3.2

, , 가

direct Euclidian

distance .

direct Euclidian distance

가

가

COD-DBSCAN .

3.2.1

가

가

[3.3]

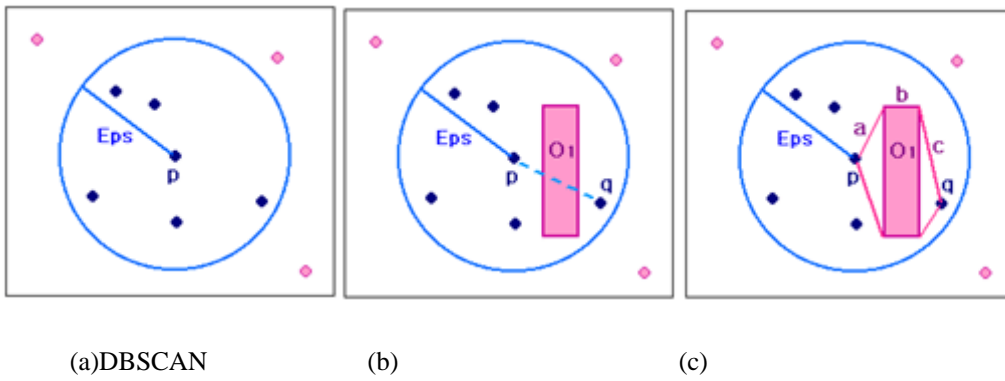
3.2.2

가

COD-DBSCAN

(Clustering with Obstructed Distance - DBSCAN)

DBSCAN



[3.3]

DBSCAN

3.3 (a) DBSCAN

p Eps MinPts

Eps

direct

Euclidian distance

(b)

p, q 가

가

p q

O_1

p q

3.3 (c) p q

p q

가 O_1

a, b, c

가 , 가

.

가 Dijkstra's [14]

.

DBSCAN , direct Euclidian distance

가 .

IV. 가

가

DBSCAN 가

가 DBSCAN-W COD-DBSCAN

가 .

가

DBSCAN , 가 가

가

DBSCAN-W .

가 ,

COD-DBSCAN .

가 .

4.1

가

가 가

4.1 .

[4.1]

	Sun Solaris 2.6
- DBMS	Informix Universal Server
	Informix Spatial Datablade Module 2.2 Informix ESQL/C Informix Datablade API GCC Compiler Microsoft VisualBasic 6.0

가

가

Sun Solaris 2.6 - DBMS Informix Universal Server . Informix Universal Server Informix-OnLine Dynamic Server DSA 2D, 3D, , ,

RDBMS

, - DBMS [15].

. Informix Spatial Datablade

Module Informix Universal Server Database

2

9 data type 45 SQL function

R-tree

[16,17,18,19,20,21].

Informix Datablade API ESQL/C

가

Microsoft VisualBasic 6.0

4.2

DBSCAN,

DBSCAN-W, COD-DBSCAN 가

가

[4.2]

DataPoint		
id	smallint (not null)	ID
location	sp2Pnt	(x,y)
region	sp2Circ	가 (x, y, r)
weight1	int	가 1
weight2	int	가 2
cluster	smallint	

[4.3]

Obstacle		
id	smallint (not null)	ID
location	sp2Pnt	(x, y)
region	sp2Box	(x1, y1, x2, y2)

DataPoint

, id가
 primary key , location 2 (x, y)
 sp2Pnt type . region
 sp2Circ type . 가
 integer type weight1 weight2 가 . cluster
 smallint type .

Obstacle

, DataPoint 가 id location . ,
 , MBR(minimum boundary rectangle)
 region
 sp2Box type .

4.3 가

4.3.1

가

가

MinPts

input parameter

Eps

DBSCAN

가

DBSCAN-W, COD-DBSCAN

DBSCAN

Eps

MinPts

가

DBSCAN-W

DBSCAN

가

가

가 가

Eps

MinPts

DBSCAN-W

COD-DBSCAN

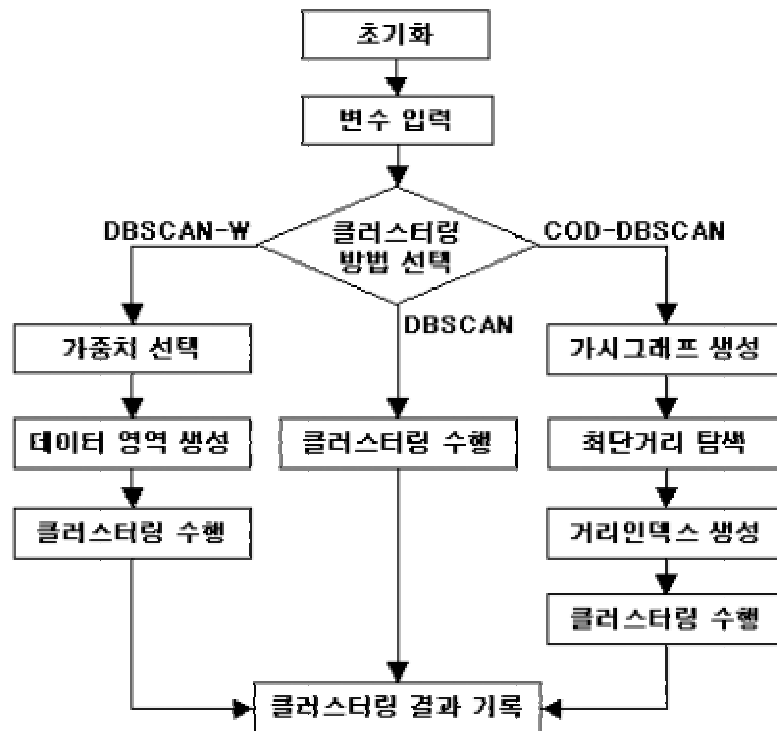
Dijkstra's

가

가

DBSCAN

4.1 가



[4.1] 가

4.3.2

가 4.4 .
 가 DBSCAN
 DBSCAN
 , Eps MinPts가 . DBSCAN-W
 가 가

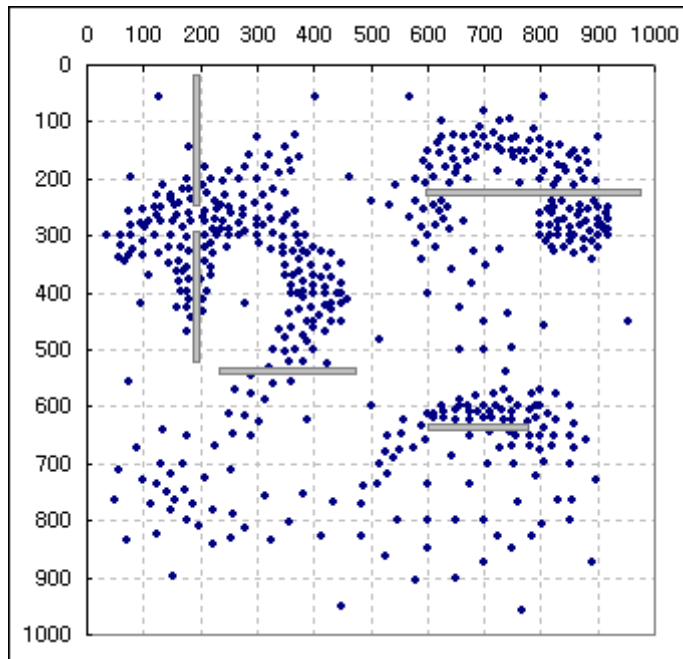
[4.4]

Eps	1~500	
MinPts	1~500	
weight	A, B	가

4.4 가

Eps MinPts 가 .
 가 (x, y) 2 (0, 0) ~
 (1000, 1000) 500 ,

5 가
 Eps 50, MinPts 5
 4.2 가 , 5

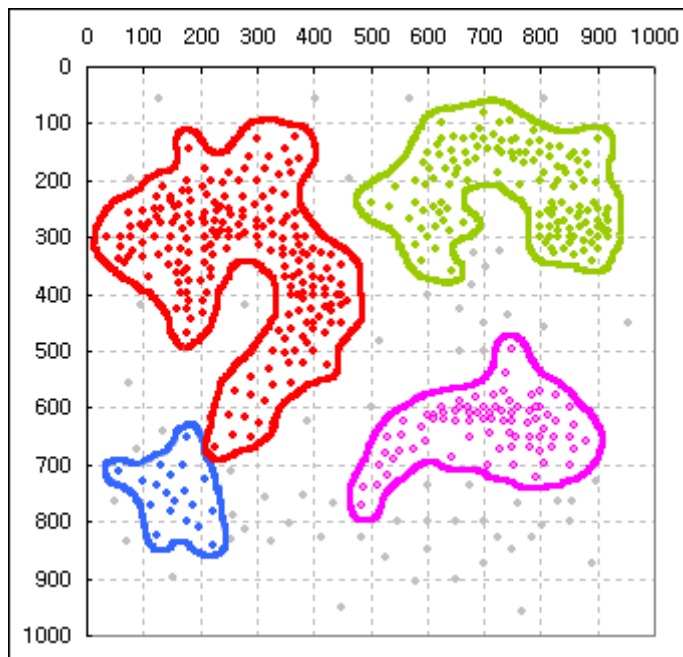


[4.2]

4.3

DBSCAN

, 4 가 . 4



[4.3] DBSCAN (Eps=50, MinPts=5)

4.4 가

DBSCAN-W

가

DBSCAN

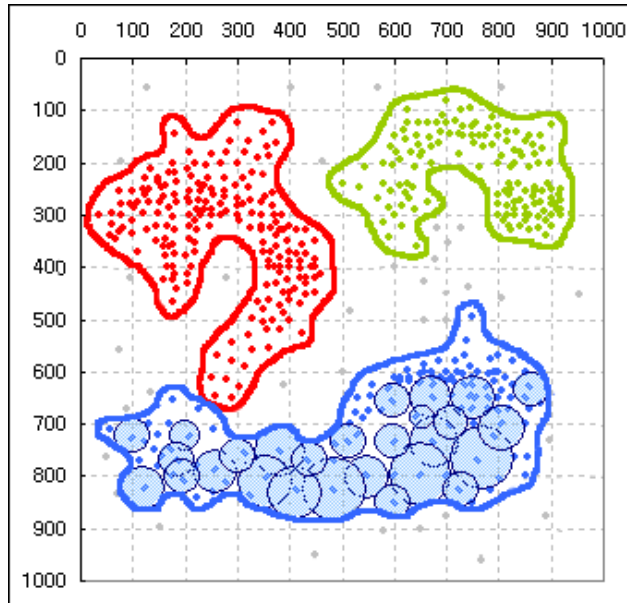
, 4.3

4.4 가

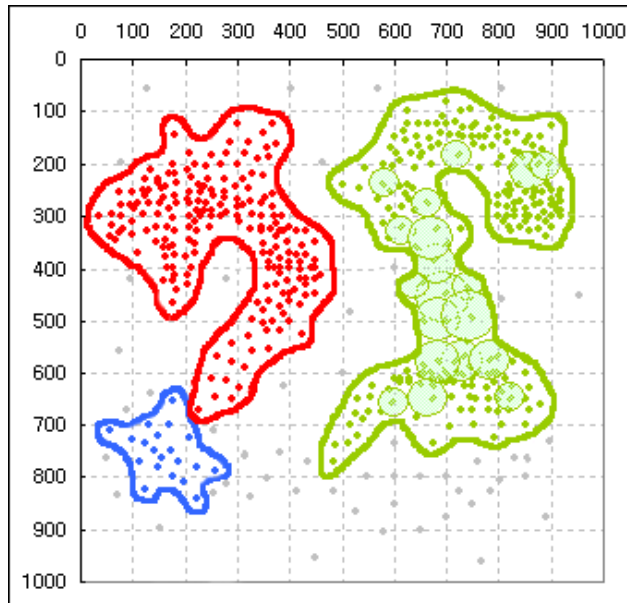
4.4 (a) weight1 가

가 가 가 ,

가



(a) weight1 가 DBSCAN-W

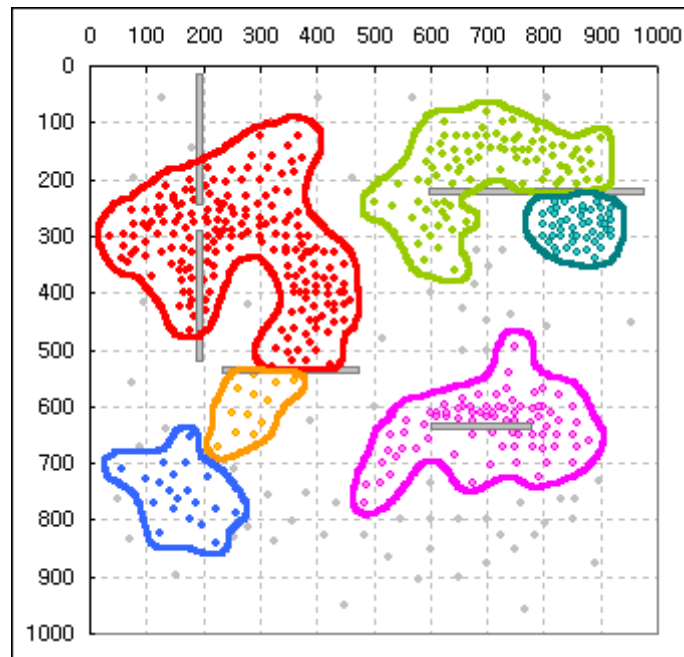


(b) weight2 가 DBSCAN-W

[4.4] DBSCAN-W (Eps=50, MinPts=5)

(b) weight2 가

가



[4.5] COD-DBSCAN (Eps=50, MinPts=5)

4.5

COD-DBSCAN

가

COD-DBSCAN

4.5

가 .

가 가

,

가

. DBSCAN COD-DBSCAN 4.3 4

가 4.5

6 .

가

,

가

가 DBSCAN-W , 가

COD-DBSCAN

IV.

가

, ,

. , ,

.

,

가

.

,

가

.

가

,

, 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, Jiawei Han, "Efficient and Effective Clustering Method for Spatial Data Mining", In Proc. of the VLDB Conference, Santiago, Chile, 20th Int, 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, Aug.1997.

- [8] Jiawei Han, Micheline Kamber, “Data Mining: Concepts and Techniques”, Morgan Kaufmann publishers,2001.
- [9] Samet, Hanan, “Spatial Data Models and Query Processing. In Modern Database Systems : The Object Model, Interoperability, and Beyond”, Addison Wesley / ACM Press, Reading, MA, 1994.
- [10] Fayad, Usama M., “Advances in Knowledge Discovery in DataBases”, AAAI Press / MIT Press, Menlo Park, CA, 1996.
- [11] Vladimir Estivill-Castro and Ickjai Lee, “AUTOCLUST : Automatic Clustering via Boundary Extraction for Mining Massive Point-Data Sets”, International Conference on Geocomputation, 2000.
- [12] M.deBerg et al, “Computational Geometry in C”, Cambridge University Press, 1994.
- [13] Michael J. Laszlo, “Computational Geometry and Computer Graphics in C++”, Prentice Hall, 1996.
- [14] , , “C ”, , pp.199-202, 1995.
- [15] Informix, Informix Universal Server Guide to SQL: Tutorial Version 9.1, Informix Press, 1997.
- [16] Informix, Informix Universal Server Guide to SQL: Reference, Informix Press, 1997.
- [17] Informix, Extending Informix Universal Server: Data Types Version 9.1, Informix Press, 1997.

- [18] Informix, Informix Spatial Datablade Module: User's Guide, Informix Press, 1997.
- [19] Informix, Informix Datablade API Programmer's Manual, Informix Press, 1997.
- [20] Informix, Developing Applications Using Informix-ESQL/C, Informix Press, 1997.
- [21] Thomas Brinkhoff, Hans Peter Kriegel, Bernhard Seeger, "Efficient Processing of Spatial Joins Using R-trees", Proceedings of the ACM, pp.237-246, 1993.
- [22] Anthony K. H. Tung, Jean Hou, Jiawei Han, "Spatial Clustering in the Presence of Obstacles", Proc. 2001 International Conference on Data Engineering, 2001.
- [23] , , , "가
", Korean Database Conference 2001 6 17 2 ,
pp.109-113, 2001.

ABSTRACT

A Density-based Spatial Clustering Algorithm Considering Weight and Obstructed Distance

*Department of Computer Science & Engineering
Ewha Institute of Science and Technology*

Lim, Hyun Sook

Spatial data mining is a process to discover interesting relationships and characteristics that exist implicitly in a spatial database. Clustering techniques work well with a large amount of data and various types of data, so that it is easy to apply clustering techniques to spatial data. The density-based clustering algorithms are effective especially at handling the locality.

DBSCAN is the locality-based algorithm, relying on a density-based notion of clustering. The density-based notion of clustering states that within each cluster, the density of the points is significantly higher than the density of points outside the cluster. The algorithm can handle the issue of noise and is successful in discovering arbitrary shaped clusters on a large amount of spatial data. But it restricts the attributes that can affect clustering results as a location and cannot manage obstacles exist in real world such as rivers, lakes and highways.

In this thesis, we present two new spatial clustering algorithms, called DBSCAN-W and COD-DBSCAN. DBSCAN-W algorithm assigns regions to objects relying on the weights and the regions have influence on density. COD-DBSCAN generates a distance index for the obstructed distance and uses it instead of the direct Euclidian distance when there exist one or more obstacles between two objects. We conduct various performance studies to show that DBSCAN-W and COD-DBSCAN produce different clustering results from DBSCAN depend on the purpose of applications.