

(Performance Comparison of Clustering Techniques for Spatio - Temporal Data)

(Nayoung Kang,), (Juyoung Kang,
) , (Hwan - Seung Yong,)

가

GPS , , .

SOM(Self-Organizing Map) ,

K - means 가 (Hierarchical Agglomerative) , , 가 가

가

가 .

ABSTRACT

With the growth in the size of datasets, data mining has recently become an important

research topic. Especially, interests about spatio-temporal data mining has been increased which is a method for analyzing massive spatio-temporal data collected from a wide variety of applications like GPS data, trajectory data of surveillance system and earth geographic data. In the former approaches, conventional clustering algorithms are applied as spatio-temporal data mining techniques without any modification.

In this paper, we focused to SOM that is the most common clustering algorithm applied to clustering analysis in data mining area, and develop the spatio-temporal data mining module based on it. In addition, we analyzed the clustering results of developed SOM module and compare them with those of K-means and Agglomerative Hierarchical algorithm in the aspects of homogeneity, separation, silhouette width and accuracy. We also developed specialized visualization module for more accurate interpretation of mining result.

: , , , SOM, 가

Keywords : Data Mining, Spatio-Temporal Data Mining, Clustering, SOM, Performance Evaluation

1.

, (natural science observation systems),

. (alphanumeric data)

[1].

[2].

가 [1]. K-means, SOM,

가 [3,4],

가 ,

가 .

SOM ,

GSTD(Generate Spatio-Temporal Data) [5]

SOM 가 가

SOM

SOM K-means ,

(Average Linkage Method), Ward 가
(homogeneity), (separation), (silhouette width),
(accuracy) 가 가 .
Insightful S-PLUS[6]

가
가 가
가 .

가

2.

가

가

가

가

MST, CLARANS, CURE

[7]가

가

[8,9]

가

가

2.1

가

[10].

가

가

. CONQUEST(CONTENT -based

Querying in Space and Time)[11], TSA-Ttree[12], Quakefinder[13]

(surveillance system)

, GPS

(trajectory)

가

GPS

(refinement)

[14],

[3,4,15,16,17].

가

K-means

(pattern discovery)

가

K-means, SOM

(agglomerative hierarchical)

가 . [16], 가
 (手信號) (haptic)
 K-means
 가 [18]. Neil Johnson
 David Hogg [3], Jonathan Owens Andrew Hunter[4] 2

SOM 가
 [3].
 , K-NN 가 [17,19,20].

(isokinetic)
 (傷害) , 가 [21],

가 .
2.2 가
 가
 , (partitioning method), (hierarchical method),
 (density-based method), (grid-based method),
 (model-based method) 가 [22].

K-means PAM, BIRCH CURE,
DBSCAN 가 [22].

SOM(Self-Organizing Map) (bayesian network)

가

[8] 가
(homogeneity), (separation), (silhouette Width),
(accuracy), (redundant Score), WAPD(
) ,
가 [9].

3. SOM

가

SOM 가 SOM
K-means, , Ward
가 4가

3.1

■ K-means

K-means K

- 1) K
- 2) K
- 3) 가 가
(Euclidean distance)
- 4)
- 5) 가 ,
2 가



bottom-up n 가
n 가 가
가 가 k 가
k

(Linkage Method) (Ward Method)
(single linkage method), (complete linkage method),
(average linkage method)

[15].

가 ,
가

S-PLUS

■ SOM(Self-Organizing Map)

SOM Kohonen

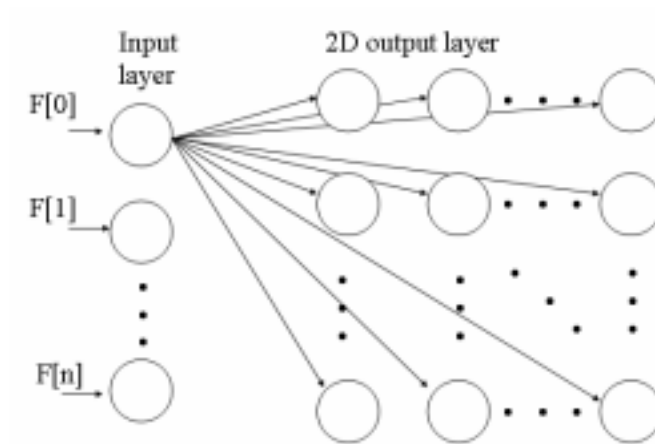
(Competitive Learning)

가 가

[23]. 1 SOM 2

(Input Layer) , 2 (Output

Layer) ,



(1) SOM

SOM

. SOM

1) K

2) $\alpha(t)$. 0 1 가 ,

3) $x(t)$

4) 가

Euclidian

$$|x(t) - m_c(t)| = \min |x(t) - m_i(t)|$$

5) 가 . 가
(Neighborhood Function)

$\lambda_{ci}(t)$ Gaussain

$$m_i(t+1) = m_i(t) + \alpha(t)\lambda_{ci}(t)[x(t) - m_i(t)]$$

$$\lambda_{ci}(t) = \exp\left(-\frac{|r_c - r_i|^2}{2\sigma^2(t)}\right)$$

6) 2 가 . SOM SOM

3.4 SOM

SOM ,
가 1

< 1> SOM

DBMS	Oracle 9i
	Oracle PRO*C Microsoft Visual C++ 6.0 Microsoft Visual Basic 6.0 ADO 2.6 Library

SOM Oracle Pro*C Oracle 9i
가 VisualBasic 6.0 ADO 2.6

3.4.1

SOM

2

< 2>

VALID	VARCHAR(20)	
ID	NUMBER(4,0)	ID
TIME	NUMBER(18,6)	가 ()
X	NUMBER(18,5)	- x
Y	NUMBER(18,5)	- y

n (frame) i n (Flow Vector)

Qi .

$$Q_i = \{f_1, f_2, f_3, \dots, f_n\}$$

4가 ..

$$f = (x, y, dx, dy)$$

x y x y , dx dy 가

x y .

dx dy .

. 1

, SOM 0 1 가 x, y 1

가 dx, dy .

dx dy 가 . dx x

. dy y

3.4.2 SOM

3

K-means

SOM .

< 3 >

ID	NUMBER(4,0)	ID
TIME	NUMBER(18,6)	가 ()
X	NUMBER(18,5)	- x
Y	NUMBER(18,5)	- y
DX	NUMBER(18,5)	가 x
DY	NUMBER(18,5)	가 y
CLUSTER	NUMBER(4,0)	

SOM

K

가

가

가

K

SOM

가

10000

3.4.4 가

가

2

가

가

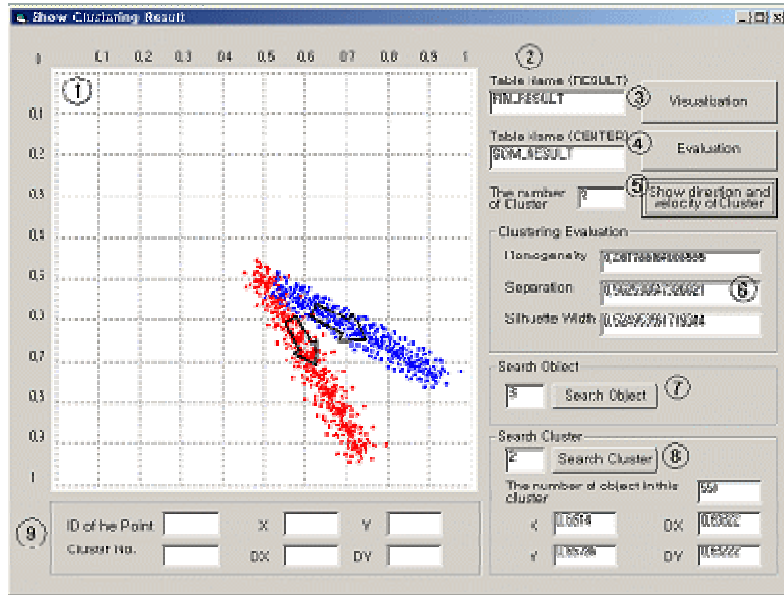
가

SOM

()

가

가



(2)

가

< 4>가

	가
	가
가	
	가
	가
,	,
ID	가
	ID, ,

3

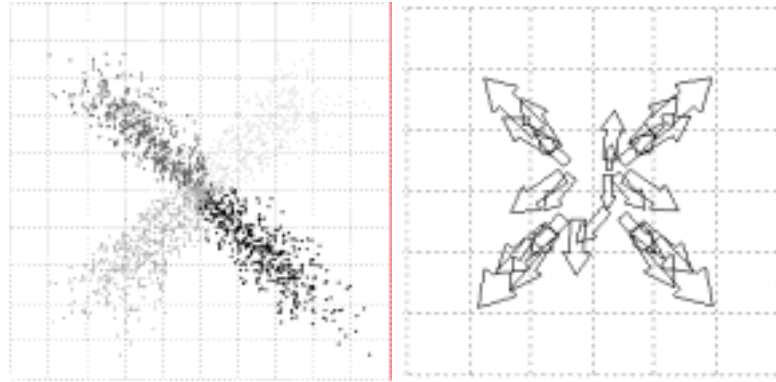
SOM

가

가

가

가



(3)

SOM

가

4. 가

4.1

2

3

가

GSTD

. GSTD

Alberta

4-a

<http://db.cs.ualberta.ca:8080/gstd>

. GSTD

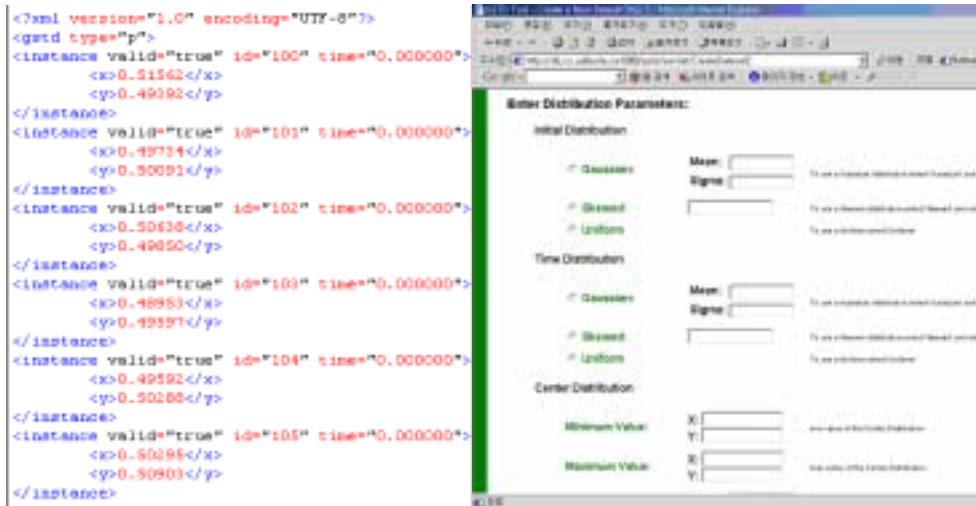
가

4-b

XML

가

[5]. GSTD



(4-a) GSTD

XML

(4-b) GSTD

GSTD가

ID Oid

(Oid, s_i, t_i)

. s_i t_i

0

. s_i t_i 0 1

가 . 5-a

,

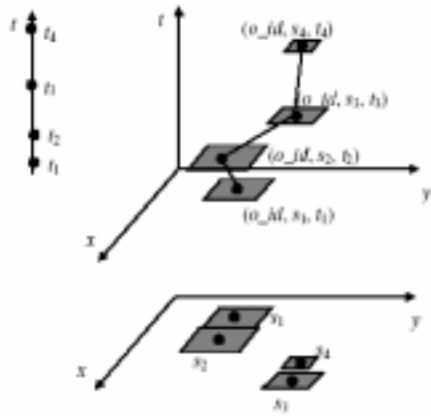
, ,

가

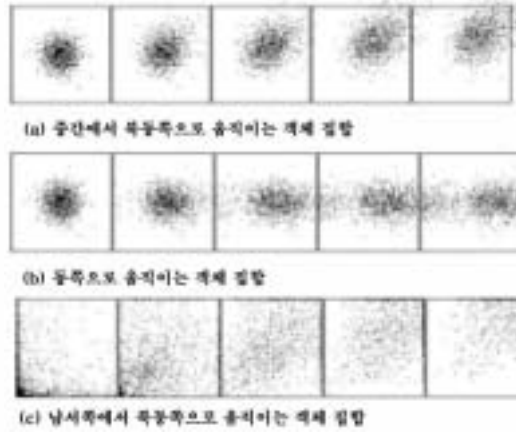
. 5-b

가

[5].



(5-a)



(4-b) GSTD

4.2

가

가

D1 ~ D6,

D7 ~ D16,

D17 가 , 17

D1 ~ D6 ,

가

1.0/s 0.1/s

D7 ~ D16

가

D17

가

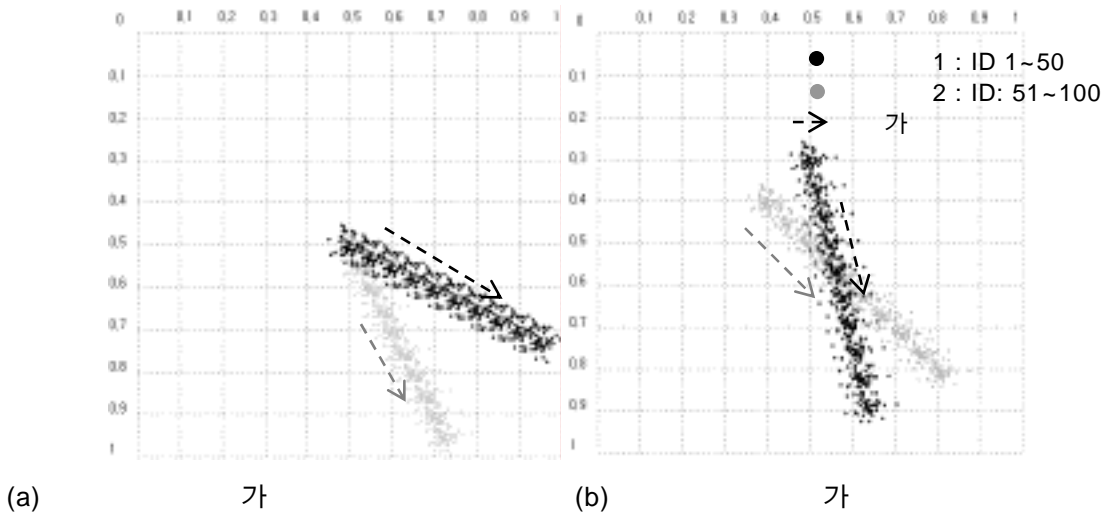
< 5 >

(D1-D17)		90	75	60	45	30	1.0/s ~ 0.1/s	(+)

		D1	D2	D3	D4	D5	D7 ~ D16	D17
		D6						

< 6 >

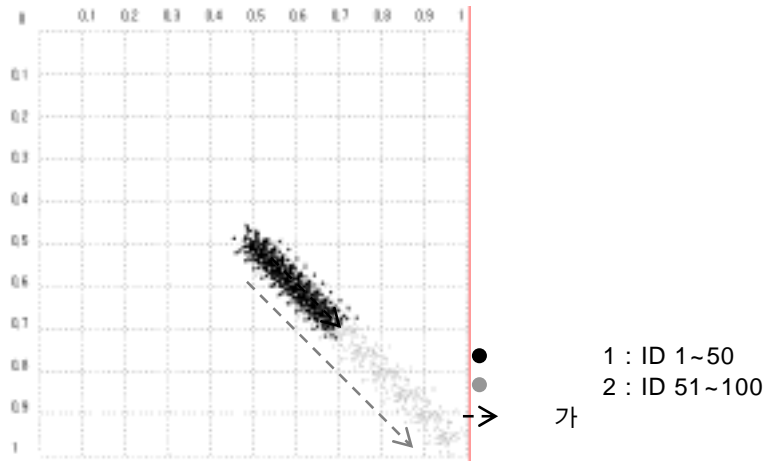
6 (a) (b) D4 D6(가 45)



(6)

7 D11(가0.6/s)

	D1~D10	D11~19	20
	2	2	4
	50	50	50
	100	100	200
	500, 500	500, 347	530, 530, 530, 530

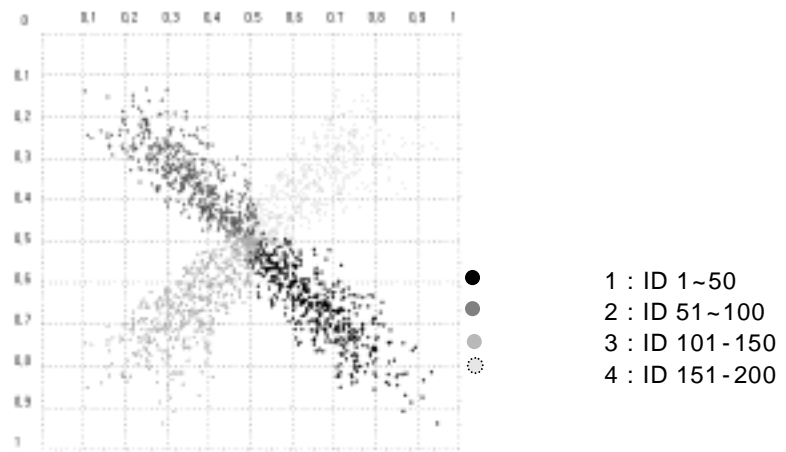


(7)

8

4

D17



(8)

4.3 가

가 (validation) ,

가 . 가

가 . 가

(Gold Standard)

가

가 [24].

가 [22].

가 가 가

■ (Homogeneity)

가

$$H_{ave} = \frac{1}{N_{point}} \sum_i D(p_i, C(p_i))$$

D , p_i i , $C(p_i)$ p_i 가, N_{point}

■ (Separation)

가

$$S_{ave} = \frac{1}{\sum_{i \neq j} N_{ci} N_{cj}} \sum_{i \neq j} N_{ci} N_{cj} D(C_i, C_j)$$

C_i C_j i j , N_{ci} N_{cj} i

j

■ (Silhouette Width)

(quality)

가 가 . 가

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$a(i)$ i

$b(i)$ i 가

■ (Accuracy)

가

$$A_{ave} = \frac{1}{N_{point}} \sum_i A_i$$

A_i i

4.3 가

가 가

가 .

D1 ~ D6

D7 ~ 16

(threshold)

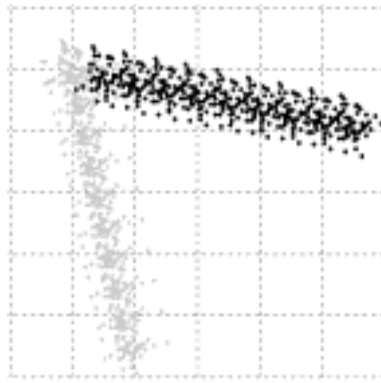
가

■ D1 ~ D6

D1 ~ D5
 가 . 9 가 75
 가 . 10, 11, 12 가
 60, 45, 30 가 60
 45 SOM Ward 가 K-means

가 30

SOM

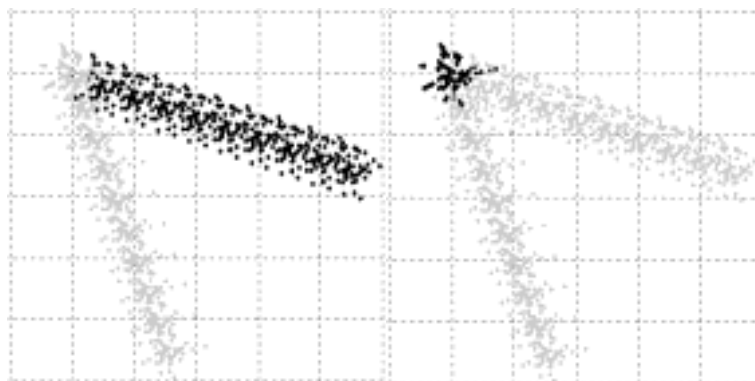


SOM, K-means, , Ward

(9)

가 75

가



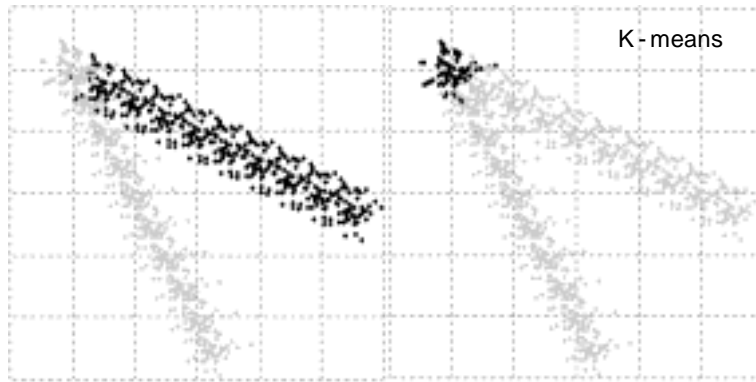
(a) SOM, Ward

(b) K-means,

(10)

가 60

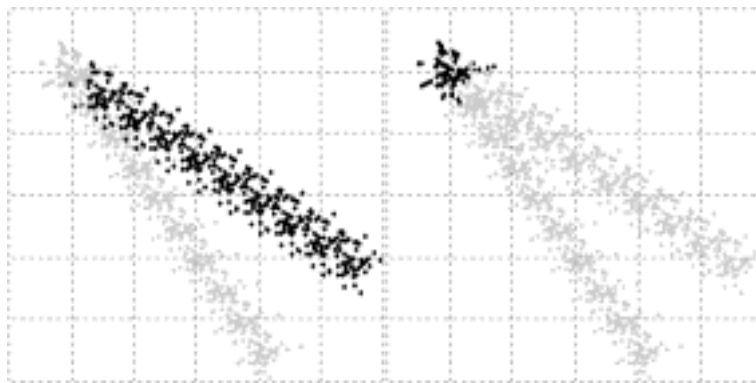
가



(a) SOM, Ward

(b) K-means,

(11) 가 45 가



(a) SOM

(b) K-means, , Ward

(4) 가 30 가

SOM

30 , K-means

75 , Ward

45

, SOM 가

SOM

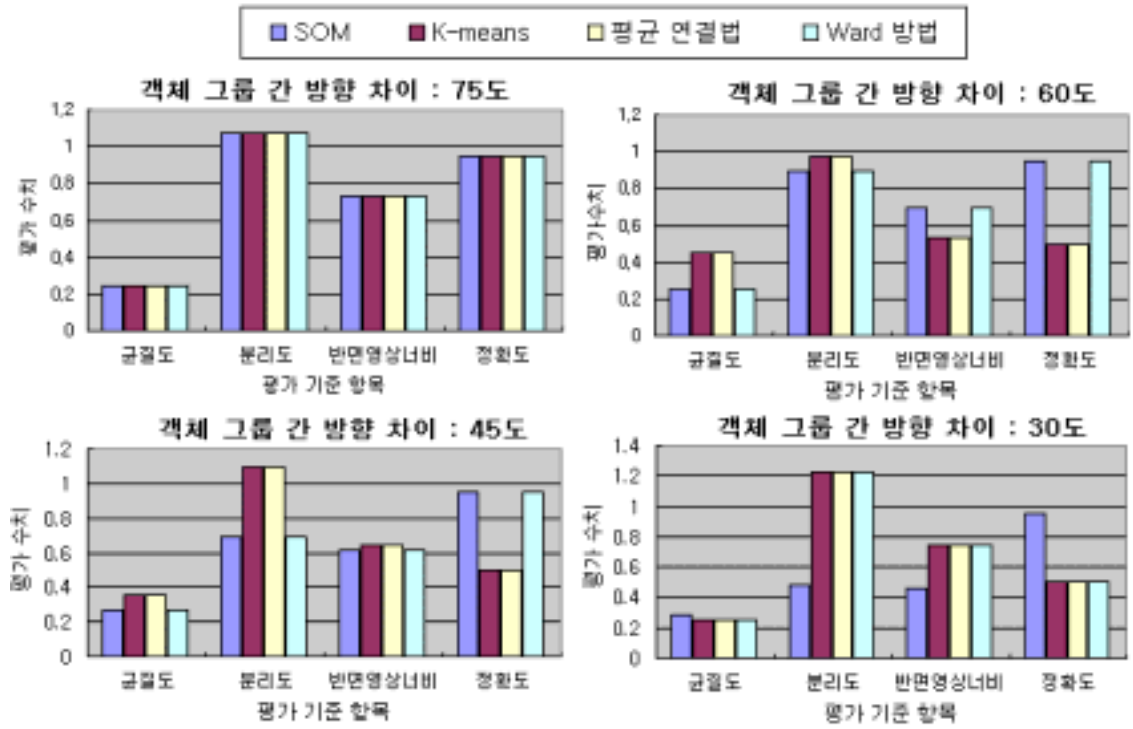
13

D1 ~ D5

가

가

가 SOM



(5) 1-5 가

■ D7 ~ D16

D7 ~ D16 가

14

가 0.6/s

가

15, 16

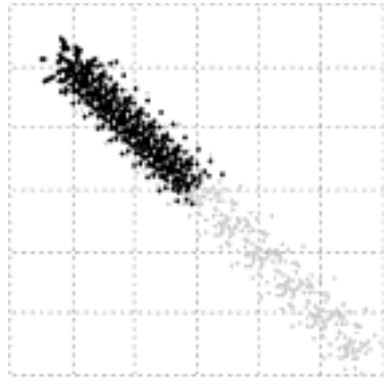
가 0.4/s, 0.3/s

가 0.4/s

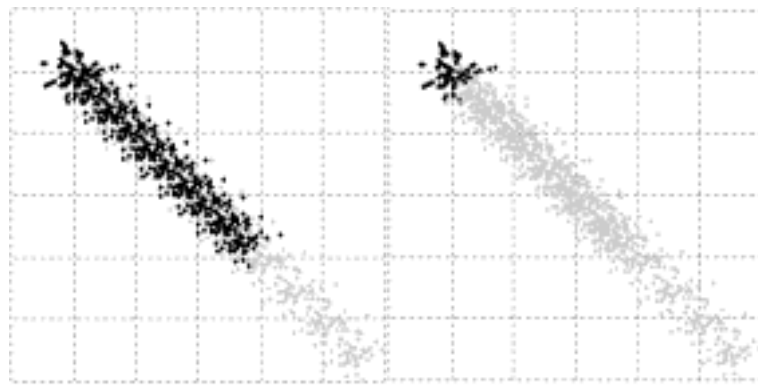
SOM Ward

가 K-means

, 가 0.3/s SOM

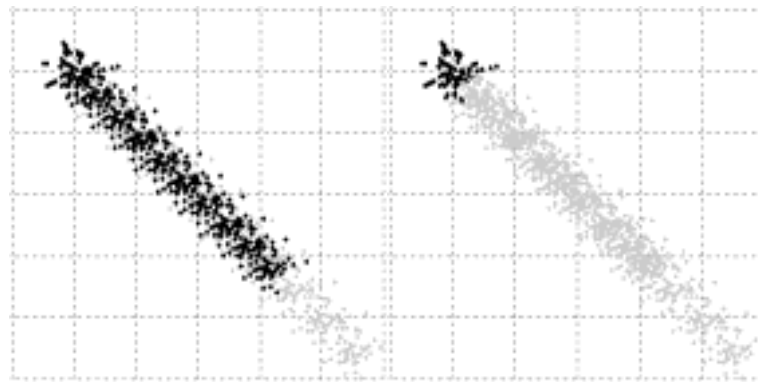


(6) 가 0.6/s 가



(a) SOM, Ward (b) K-means,

(7) 가 0.4/s 가



(a) SOM (b) K-means, , Ward

(8) 가 0.3/s 가

Ward

0.3/s, K-means

0.4/s

SOM

가

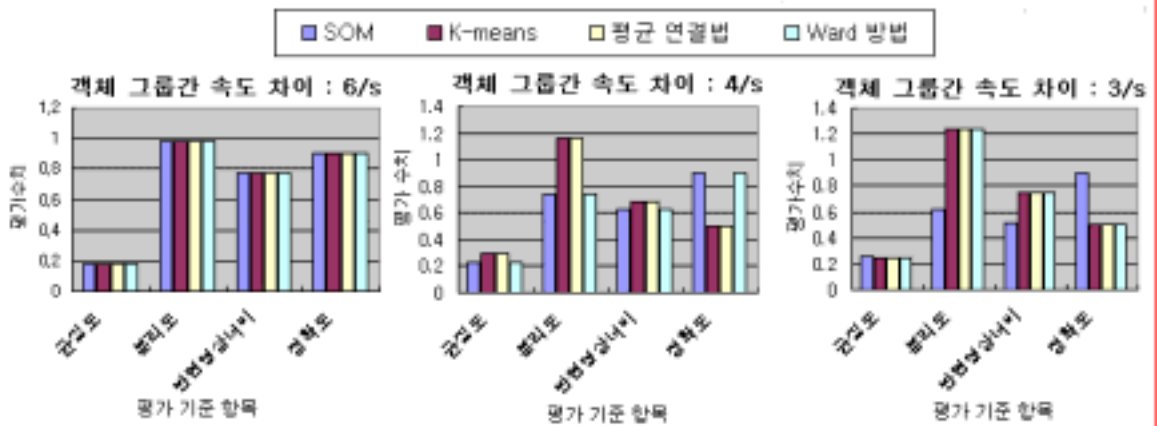
SOM

17

11, 13, 14

가

SOM



(9)

11, 13, 14

가

18

4

D17

7

가

가

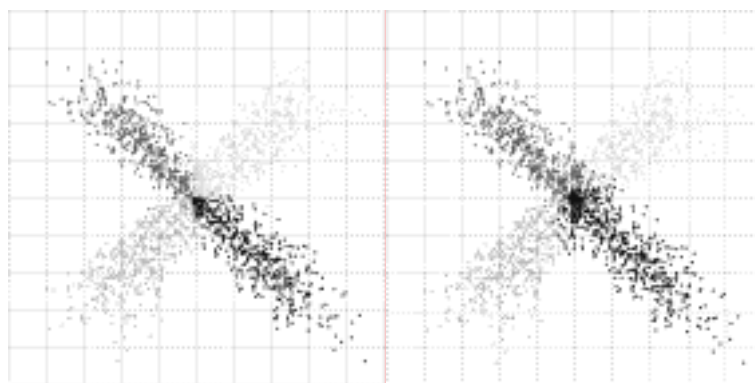
SOM, K-means, Ward

가

18

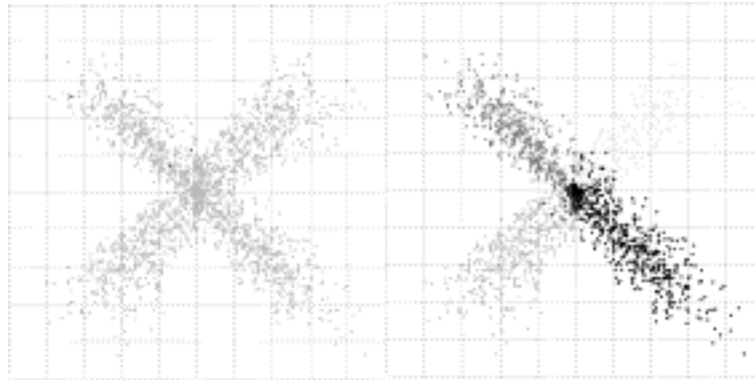
SOM

가



(a) SOM

(b) K-means

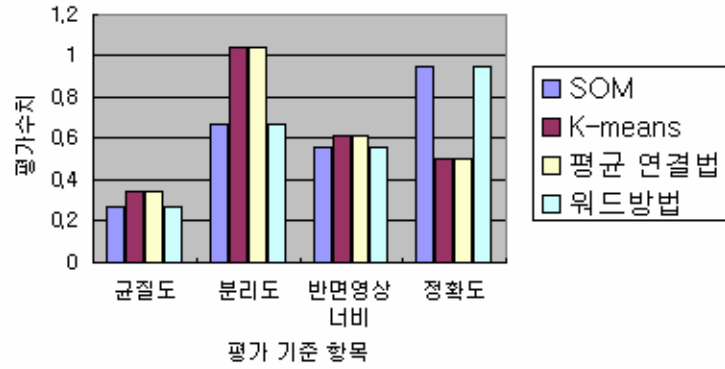


(c)

(d) Ward

(10) 17 가

데이터세트 D17에 대한 클러스터링 결과



11 D17 가

가

가

가

SOM

가

가

SOM

가

SOM

(Classification)

가

가 30

가

0.3/s

SOM

K-means

SOM 가

K-means가 SOM

SOM

SOM

SOM

K-means,

가

3가

17가

가

가

가 가

가

K-means

SOM

가

SOM

가

SOM

가

가

가

가

- [1] J.F. Roddick, and M. Spiliopoulou, "A bibliography of Temporal, Spatial and Spatio-temporal Data Mining Research", Explorations 1(1), pp.34-38, 1999
- [2] J.F. Roddick, and B. G. Lees, "Paradigms for Spatial and Spatio-Temporal Data Mining, Geographic Data Mining and Knowledge Discovery", Miller, H & J. Han (eds). Taylor & Francis, pp.33-50, London, 2001.
- [3] N. Johnson and D Hogg, "Learning the Distribution of Object Trajectories for Event Recognition", Image and Vision Computing, 14(8), pp.609-615, 1996
- [4] J. Owens and A. Hunter, "Application of the Self-Organizing Map to Trajectory Classification", Third IEEE International Workshop on Visual Surveillance (VS'2000), pp.77-84, 2000
- [5] Y. Theodoridis, J. R.O. Silva, and Mario A. Nascimento, "On the Generation of Spatiotemporal Datasets", In Proc. of the 6th Int'l Symposium on Large Spatial Database(SSD), pp.147-164, 1999
- [6] <http://www.insightful.com/>
- [7] , , " ", 2 , 1999
- [8] Chen G, Jaradat SA, Banerjee N, Tanaka TS, Ko MSH, Zhang MQ. "Evaluation and comparison of clustering algorithms in analyzing ES cell gene expression data", Statistica Sinica, pp.241-262, 2000
- [9] Yeung, Haynor, Ruzzo: Validating Clustering for Gene Expression Data. Technical Report UW-CSE-00-01-01, 2000
- [10] M. Erwig, R. H. Gutting, M. Schneider, M. Vazirgiannis, "Spatio-Temporal Data Types: An Approach to Modeling and Querying Moving Objects in Databases",

Technical Report, CHOROCHRONOSTR -97-08, 1997

[11] P. Stolortz, H. Nakamura, E. Mesrobian, and et al., "Fast Spatio-Temporal Data Mining of Large Geophysical Datasets", In Proc. of the First International Conference on Knowledge Discovery and Data Mining, AAAI Press, pp.300-305, 1995

[12] U. Fayyad, D. Haussler, and P. Stolorz, "KDD for science data analysis: Issues and examples", In Proc. of the Second International Conference on Knowledge Discovery and Data Mining (KDD-96), pp.50-56. CA: AAAI Press, 1996.

[13] C. Shahabi, X. Tian, and W. Zhao. "TSA-tree: A wavelet-based approach to improve the efficiency of multi-level surprise and trend queries", In The 12th International Conference on Scientific and Statistical Database Management, pp.55-68 , SSDBM, 2000

[14] S. Rogers, P. Langley, and C. Wilson, "Mining GPS data to augment road models", In Proc. of the 5th International Conference on Knowledge Discovery and Data Mining, ACM Press, pp.104-113, San Diego CA, 1999

[15] S. Gafney and P. Smyth, "Trajectory clustering with mixtures of regression models", In Proc. of the 5th International Conference on Knowledge Discovery and Data Mining, ACM Press, pp.104-113, San Diego CA, 1999

[16] C. Stauffer and W. Eric L. Grimson, "Learning patterns of activity using real-time tracking", IEEE Trans. PAMI, vol. 22, pp.747-757, 2000

[17] N. Sumpter and A. Bulpitt, "Learning spatio-temporal patterns for predicting object behaviour", Technical report, University of Leeds, School of Computer Studies, The University of Leeds, UK. 1998

[18] J. Eisenstein, S. Ghandeharizadeh, L. Huang, C. Shahabi, G. Shanbhag and R. Zimmermann, "Analysis of Clustering Techniques to Detect Hand Signs", Int'l

Symposium on Intelligent Multimedia, Video and Speech Processing, Hong Kong, 2001

[19] P. Remagnino, T. Tan, and K. Baker. "Agent orientated annotation in model based visual surveillance", In ICCV, pp 857-862, 1998

[20] S. Handley, P. Langley and F. Rauscher, "Learning to predict the duration of an automobile trip", In Proc. of the 4th International Conference on Knowledge Discovery and Data Mining, AAAI Press, pp.219-223, New York, 1998

[21] P.C. Juan and L.C. Ignacio, Discovering Similar Patterns in Time Series, In Proc. In Proc. of the 4th International Conference on Knowledge Discovery and Data Mining, AAAI Press, pp.497-505, New York, 1998

[22] J. Han and M. Kamber, "Data Mining: Concepts and Techniques", Morgan Kaufmann, 2000

[23] , (), , pp.169-189, 2001

[24] M. Halkidi, Y. Batistakis and M. Vazirgiannis, "On Clustering Validation Techniques", Journal of Intelligent Information Systems, pp.107-145, 2001