



:



:

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

(... .)

(.)

.

:

// :

:

.

.

.

.....

:

// :

:

:

را به من آموخت.

هستند.

برایم

% ,

:

.....
.....
.....
.....
.....

:

.....
.....
.....

:

.....
.....
.....
.....
.....

.....AMM

.....
.....
.....
.....
.....
.....
.....
.....
.....
.....

:

.....

.....

.....

.....

.....

.....

:

.....

.....

..... ()

..... ()

..... ()

..... C_i ()

..... *NMI* C_i ()

..... (b) (a) ()

..... (b) (a) C_1 *AMM* ()

C_1 P^a (a) *AMM* ()

..... P^a C_1 P^{b*} (b) ()

..... P^b P^a C_1 *AMM* ()

..... (a) ()

..... P^{b*} (b) P^a ()

..... ()

..... ()

..... ()

..... ()

..... ()

(B (A *NMI* C_i ()

..... (C

..... ()

..... ()

..... ()

..... ()

.....

()

.....

()

.....

()

.....

()

.....

()

•
•

.()

: .()

.()

()

-
- 1 Jain
 - 2 Faceli
 - 3 Strehl and Gosh
 - 4 Data Mining
 - 5 Classification
 - 6 Classifier
 - 7 Unsupervised
 - 8 Tracking

)

)

.(

.(b

b

a

.(

)

-
- ¹ SubClass
 - ² Fred and Jain
 - ³ Parvin
 - ⁴ Robustness
 - ⁵ Novelty
 - ⁶ Stability
 - ⁷ Flexibility
 - ⁸ Topchy
 - ⁹ Fred and Lourenco
 - ¹⁰ Ayad and Kamel

)

:(

()

)

)

(

K-means

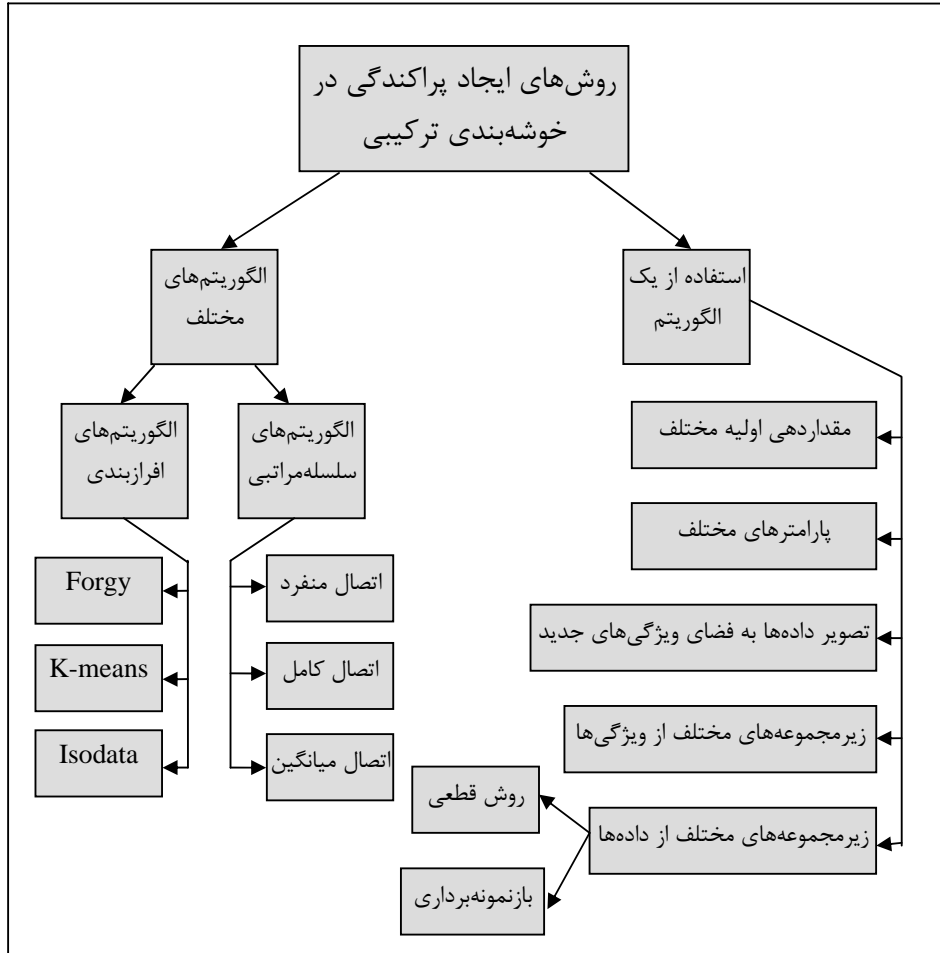
(

-
- ¹ diversity
 - ² Consensus Function
 - ³ Hierarchical
 - ⁴ Duda
 - ⁵ Partitional
 - ⁶ Jain and Dubes
 - ⁷ Kaufman and Rosseeuw
 - ⁸ Man and Goth

(
b a
)

()
()
()
()
()
()
()
()

¹ Initialization
² Barthelemy and Leclerc
³ Features
⁴ Fern and Brodley
⁵ Resampling
⁶ Dudoit and Fridlyand



()

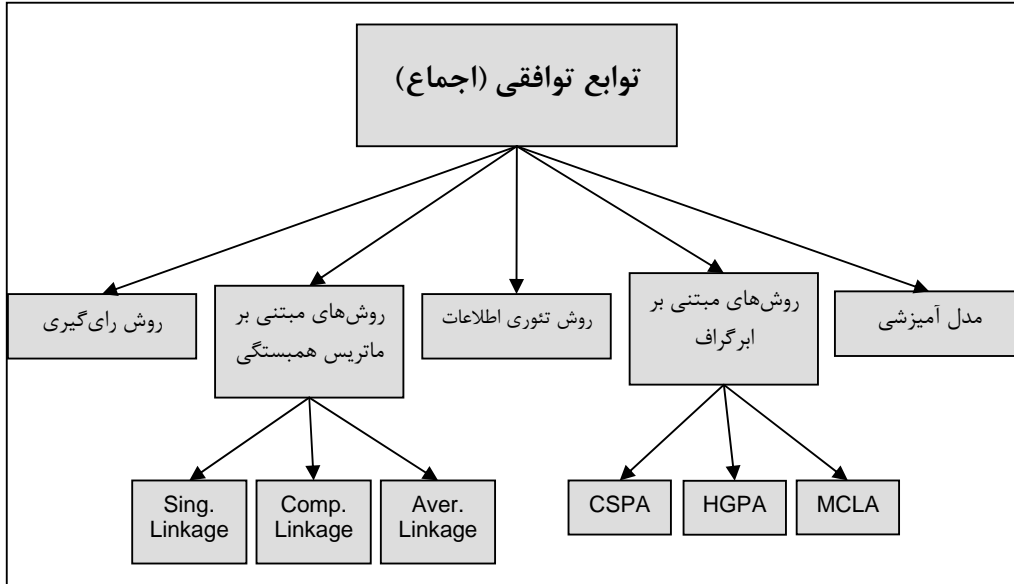
()

()

(EAC¹)

¹ Evidence Accumulation Clustering

() . ()



()

¹ Supervisor

² Train

)

.(

)

.(

"

"

()

()

)

()

.(

¹ Fischer and Buhmann

² Kuncheva and Whitaker

³ Kuncheva and Hadjitodorov

⁴ Fern and Lin

)

.(

()

.()

d

()

¹ Partitions

² Robust

³ Pairwise

()

()

¹ Classification

² Easy

³ Intermediate

⁴ Hard

⁵ Class

⁶ Experimental Results

) ()

() .(

)

.(

)

()

(b a)

(

)

() .(

.()

()

¹ Law

² Shamiry and Tishby

³ Lange

⁴ Breckenridge

⁵ Fridlyand and Dudoit

⁶ Levine and Domany

⁷ Rakhlin and Caponnetto

⁸ Roth and Lange

⁹ Cluster Validity

¹⁰ Ben-Hur

()

()

(NNR⁶)

)

()

.(

:

.()

()

()

¹ Jaccard Coefficient

² Estivill-Castro and Yang

³ Support Vector Machine

⁴ Outliers

⁵ Moller and Radke

⁶ Nearest Neighbor Resampling

⁷ Brandsma and Buishand

⁸ Inokuchi

⁹ Kernelized Validity Measure

¹⁰ Xie - Beni

¹¹ Das and Sil

¹² Full Ensemble

:

() (SNMI¹)

) (NMI²)

()

()

()

()

d

()

()

¹ Sum of Normalized Mutual Information

² Normalized Mutual Information

³ Multiobjective

⁴ Brossier

⁵ Ultrametric Distance Matrices

⁶ Lapointe and Legendre

$x_j \quad x_i \quad d_c(i,j)$
() .()

()
)
) ()
" ()
"

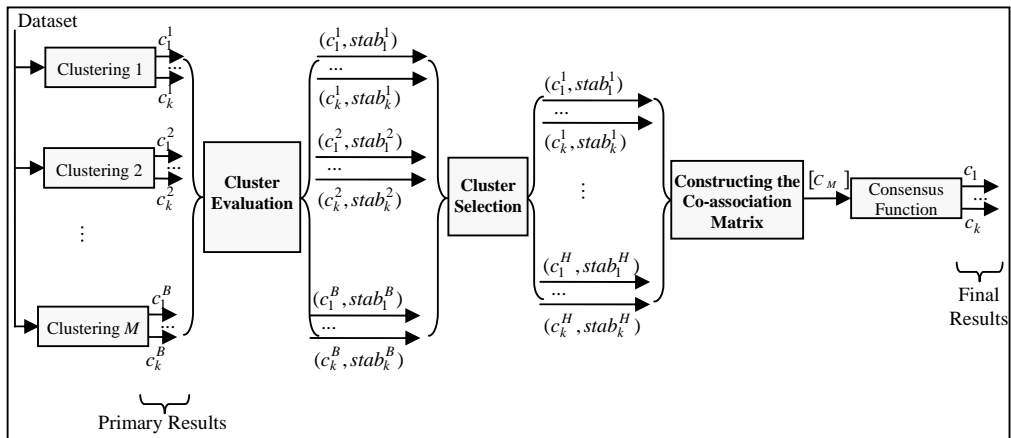
()

¹ Associated Cophenetic Matrix

² Banfield

³ Minimum Spanning Tree

•
•



()

K-means

K-means

.()

.

$$D \quad C_i \quad g_j(C_i, D)$$

:()

$$A_j \quad f_j \quad \bullet$$

$$C_i \quad g_j(C_i, D)$$

$$. \quad (\quad) \quad A_j \quad f_j$$

$$g_j(C_i, D) >$$

$$. \quad f_l \quad f_j \quad C_i \quad g_l(C_i, D)$$

$$l \quad C_i \quad j \quad C_i$$

$$f_j \quad C_l \quad C_i \quad g_j(C_i, D) = g_j(C_l, D)$$

j

$$. \quad (\quad)$$

¹ Goodness

² Perturbation

(NMI)

.()

C_i

"

$(P(D))$

(C_i)

"

$sim(C_i, P(D))$

$sim(C_i, P(D))$

$g_i(C_i, D)$

D

P

C_i

$P(D)$

C_i

$sim(C_i, P(D))$

C_i

D

D/C_i

$P(D)$

$.P_1 = \{C_i, D/C_i\}$

P_1

$D/C_i C_i$

```

For  $l:=1$  to  $M$  do
    Resample  $D$  to obtain the perturbed data set  $D'$ ;
    Run K-means over  $D'$  to obtain  $P(D')$ ;
    Re-labeling  $P(D')$  to  $P(D)$ ;
    Compute  $score[l] = sim(C_i, P(D))$ ;
End
 $g_j(C_i, D) := average$  of  $score[l]$ ;
    
```

C_i ()

.

C^* . D/C^* C^* $P(D)$

C_i %

$P_2 = \{C^*, D/C^*\}$ P_2 . D/C^*

) () (NMI)

(

(MI) . $P_2 P_1$

$P_2 P_1$ NMI . NMI

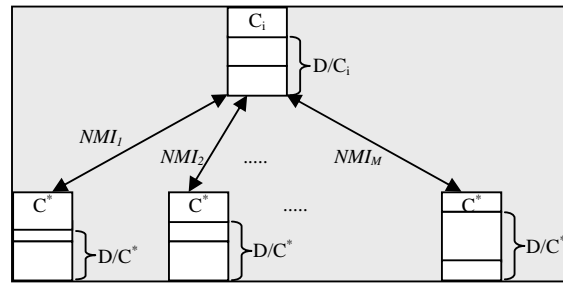
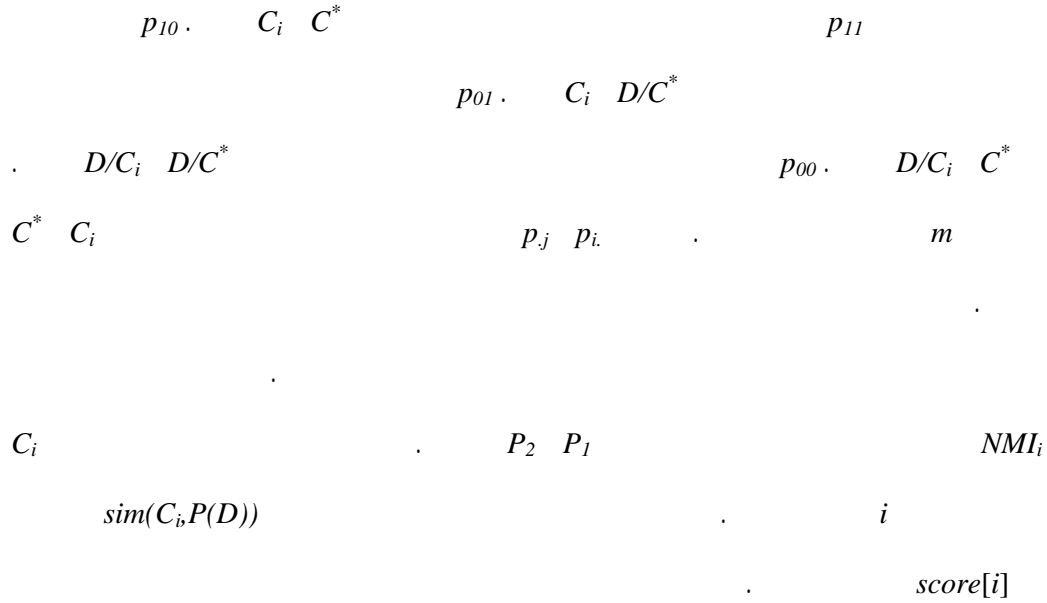
$$NMI(P_1, P_2) = \frac{MI(P_1, P_2)}{\frac{-1}{2m} \left(\sum_{i=0}^1 p_i \log \frac{p_i}{m} + \sum_{j=0}^1 p_j \log \frac{p_j}{m} \right)}$$

$$MI(P_1, P_2)$$

¹ Mutual Information (MI)

$$MI(P_1, P_2) = \sum_{i=0}^1 \sum_{j=0}^1 \frac{r_{ij}}{m^2} \log \frac{mp_{ij}}{r_{ij}} \quad ()$$

$$r_{ij} = P_i.P_j, \quad P_i = P_{i0} + P_{i1}, \quad P_j = P_{0j} + P_{1j}$$



NMI C_i ()

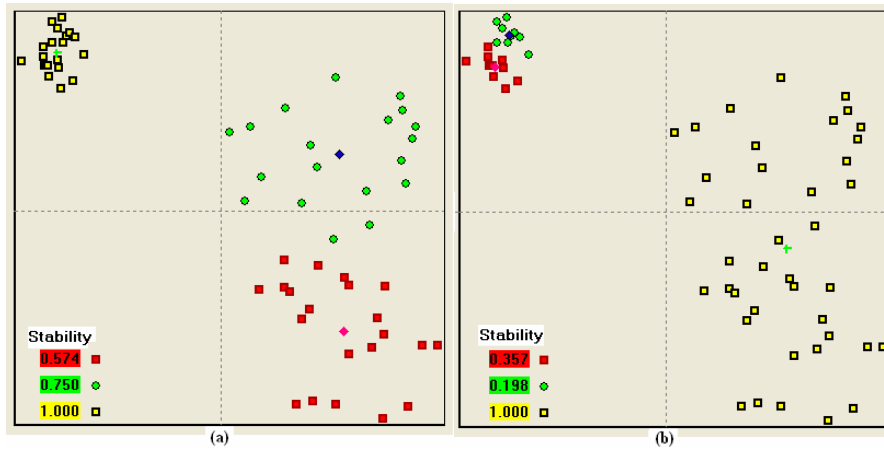
$$Stability(C_i) = \frac{1}{M} \sum_{i=1}^M NMI_i \quad ()$$

M

K-means

(b)

(a)



(b)

(a)

()

%

)

(

(b)

%

¹ Reference set

(%)

(%)

(%)

C^*

C^*

:

C^*

(C_i %)

C^*

C^*

$n_{merge(i)}$

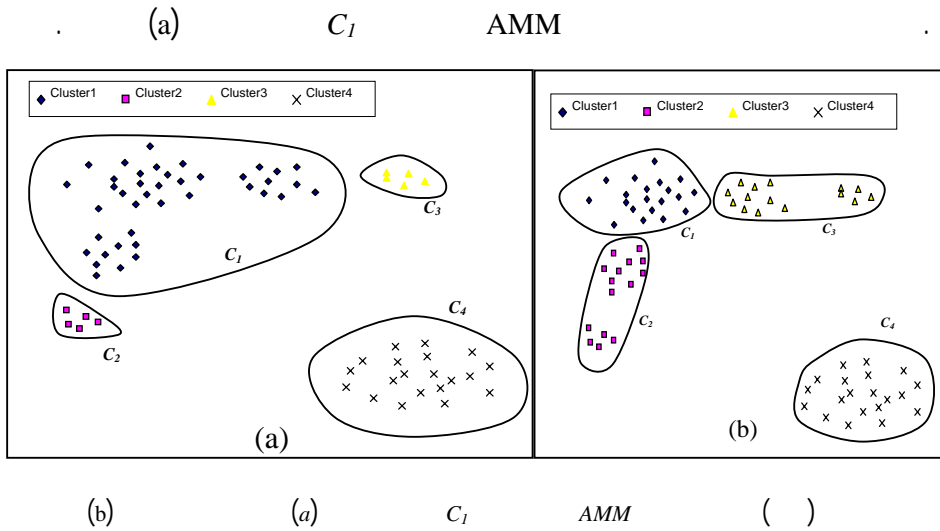
C^*

$$Stability(C_i) = \frac{1}{M} \sum_{i=1}^M w_i NMI_i \quad ()$$

$$w_i = \frac{1}{n_{merge(i)}} \quad ()$$

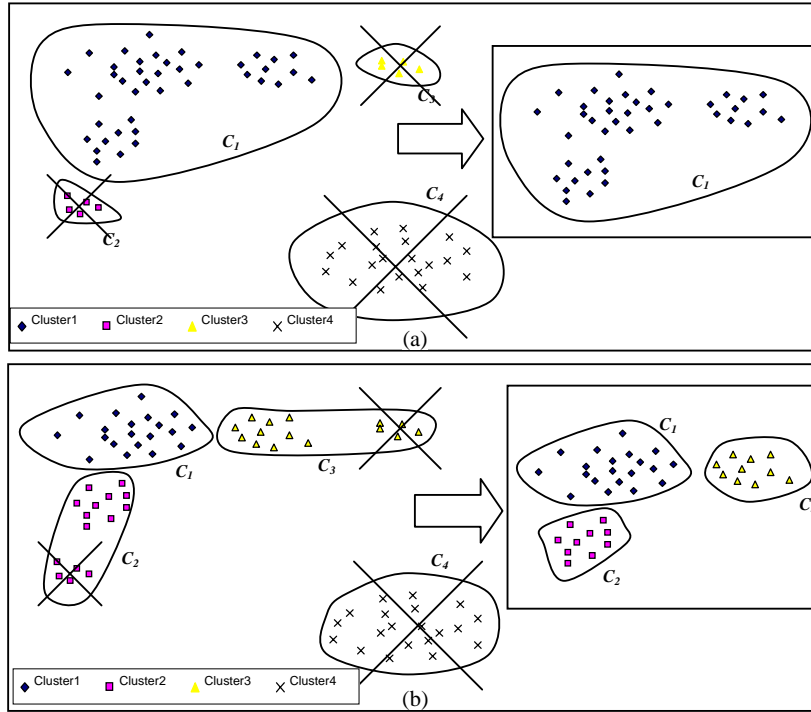
(AMM)

AMM



¹ Alizadeh-Moshki-Minaei (AMM)

² Edited Normalized Mutual Information (ENMI)



P^{b*} (b) C_1 P^a (a) .AMM ()
 P^a C_1

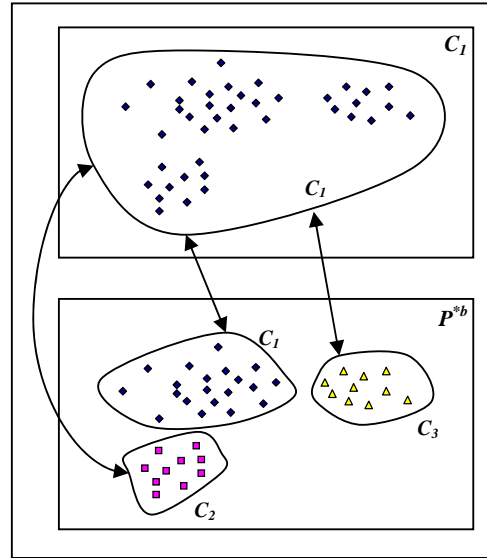
C_1 ((a)) P^a

((b)) P^b

C_1

P^b

P^b P^a



P^b P^a C_1 AMM ()

$$NMI(P^a, P^b) = \frac{-2 \sum_{i=1}^{k_a} \sum_{j=1}^{k_b} n_{ij}^{ab} \log \left(\frac{n_{ij}^{ab} \cdot n}{n_i^a \cdot n_j^b} \right)}{\sum_{i=1}^{k_a} n_i^a \log \left(\frac{n_i^a}{n} \right) + \sum_{j=1}^{k_b} n_j^b \log \left(\frac{n_j^b}{n} \right)} \quad ()$$

$C_j^b \in P^b$ $C_i^a \in P^a$ n_{ij}^{ab} n

()

P^a C_i

$$\begin{aligned}
 & C_i \\
 & n_{ij}^{ab} \\
 & (\quad) \quad n_j^b \quad C_j^b \in P^b \quad C_i^a \in P^a \\
 &) \quad P^a \quad k_a \\
 & P^b \quad P^a \quad C_i \quad \text{AMM} \quad (C_i
 \end{aligned}$$

$$\text{AMM}(C_i, P^{b*}) = \frac{-2 \sum_{i=1}^i \sum_{j=1}^{k_{b*}} n_j^{b*} \log \left(\frac{n_j^{b*} \cdot n}{n_i^a \cdot n_j^{b*}} \right)}{\sum_{i=1}^i n_i^a \log \left(\frac{n_i^a}{n} \right) + \sum_{j=1}^{k_{b*}} n_j^{b*} \log \left(\frac{n_j^{b*}}{n} \right)} \quad ()$$

$$\begin{aligned}
 & \log \quad n_j^{b*} \text{ سازی} \quad P^b \quad C_i \quad P^{b*} \\
 & (\quad C_i \quad) \quad i
 \end{aligned}$$

$$\text{AMM}(C_i, P^{b*}) = \frac{-2 \sum_{j=1}^{k_{b*}} n_j^{b*} \log \left(\frac{n}{n_i^a} \right)}{n_i^a \log \left(\frac{n_i^a}{n} \right) + \sum_{j=1}^{k_{b*}} n_j^{b*} \log \left(\frac{n_j^{b*}}{n} \right)} \quad ()$$

$$\begin{aligned}
 & j \quad \log \\
 & \text{AMM}
 \end{aligned}$$

$$\text{AMM}(C_i^a, P^{b*}) = \frac{-2 \log \left(\frac{n}{n_i^a} \right) \sum_{j=1}^{k_{b*}} n_j^{b*}}{n_i^a \log \left(\frac{n_i^a}{n} \right) + \sum_{j=1}^{k_{b*}} n_j^{b*} \log \left(\frac{n_j^{b*}}{n} \right)} \quad ()$$

$$C_i^a \in P^a \quad n_{ij}^{ab} \quad n$$

$$AAMM(C_i) = \frac{1}{M} \sum_{j=1}^M AMM(C_i, P_j^{b*}) \quad (1)$$

AMM

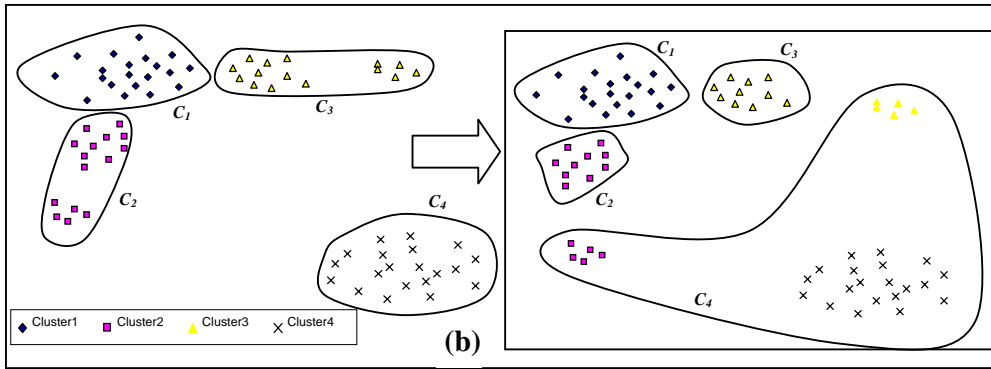
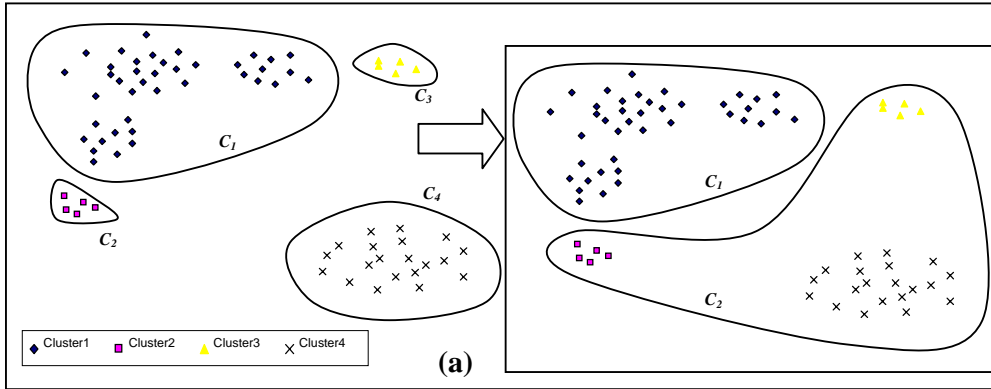
AMM

C_i

C_i

¹ Average AMM

² Edited NMI

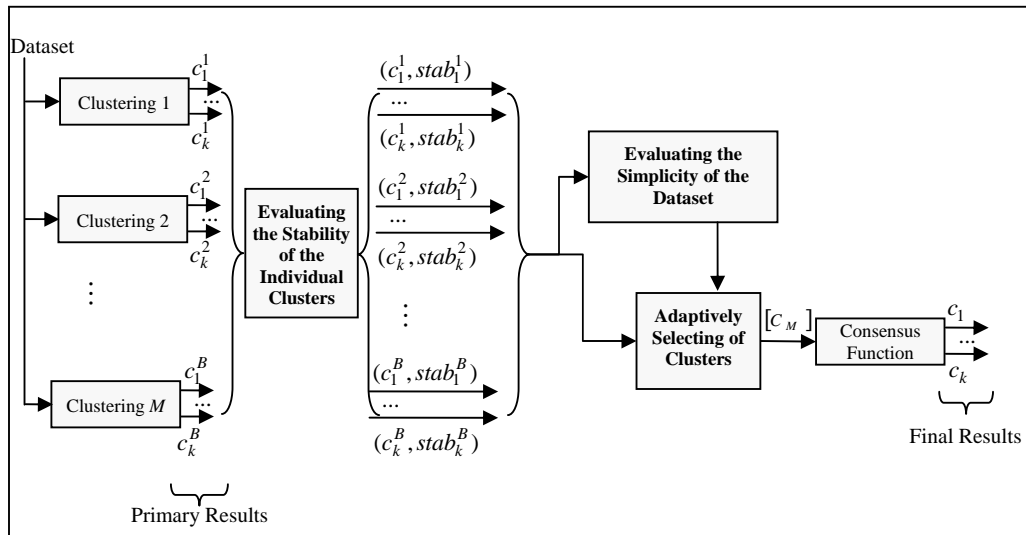


(b) P^a

(a) .

P^{b^*}

()



()



$$Stability(P) = \frac{1}{N} \sum_{i=1}^k |C_i| Stability(C_i) \quad ()$$

k

N

$|C_i|$

$$Simplicity(D) = \frac{1}{B} \sum_{i=1}^B Stability(P_i) \quad ()$$

[, ,]

M

:

%

:

•

• : %

• : %

□

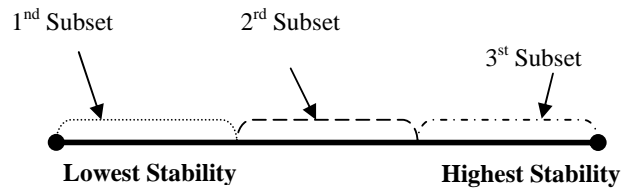
f

$f(\text{Simplicity of Dataset}) = \text{Special Subset of Clusters}$ ()

)

(

()



()

AMM

□

)

(

%

```
S := {};           //Subset of Selected Clusters
T           //Sorted Total Primary Clusters
Initialization
  T := Sort total clusters according to AMM;
  S := {the most stable cluster};
For i := 2 to length (T)
  current := i-th cluster of T;
  similar := Find the most similar cluster in S with the current
  If distance(current,similar)>th
    Add current into S;
  End If;
End For;
Return S;
```

()



```
S := {};           //Subset of Selected Clusters
T           //Sorted Total Primary Clusters
Initialization
  T := Sort total clusters according to AMM;
  S := {the most stable cluster};
For i := 2 to length (T)
  current := i-th cluster of T;
  For j := 1 to length(S)
    temp := j-th cluster of S
    If similarity(current,temp)>th1
      Add temp into similar;
    End If;
  End For;
  If mean_distance(current,similar)>th2
    Add current into S;
  End If;
End For;
Return S;
```

()

S



%

d

n

$n \times d$

) K-means

Iris

(k=3

×

```
S := {}; //Subset of Selected Clusters
G := Apply a clustering technique over all primary clusters
For i := 1 to length(G)
  current := i-th group of G;
  If size(current)>1
    Add the most stable cluster of in the current into S;
  End If;
End For;
Return S;
```

()

(EAC)

m (EAC)

$n \times n$

$$C(i, j) = \frac{n_{i,j}}{m_{i,j}}$$

()

$m_{i,j}$

$j \ i$

$n_{i,j}$

EAC

(EEAC)

EEAC

$$C(i, j) = \frac{n_{i,j}}{\max(n_i, n_j)}$$

()

n_j

i

n_i

$n_{i,j}$

j

$j \ i$

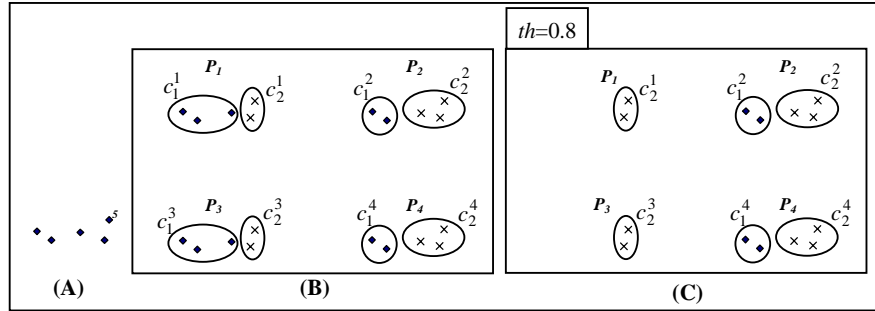
$n_j \ n_i$

$$n_i, n_j \leq B \leq k \times B$$

(A)

$P_4 \quad P_1$

.(B))



(B)

(A.NMI

C_i

()

(C)

:

$$Stability(c_2^1) = Stability(c_2^3) = 1$$

$$Stability(c_1^2) = Stability(c_1^4) = 1$$

$$Stability(c_2^2) = Stability(c_2^4) = 0.82$$

$$Stability(c_1^1) = Stability(c_1^3) = 0.55$$

.(C)

$$C(1,2) = \frac{2}{\max(2,2)} = \frac{2}{2} = 1$$

$$C(1,3) = C(2,3) = \frac{0}{\max(2,2)} = \frac{0}{2} = 0$$

$$C(3,4) = C(3,5) = \frac{2}{\max(2,4)} = \frac{2}{4} = 0.5$$

$$C(4,5) = \frac{4}{\max(4,4)} = \frac{4}{4} = 1$$

:

$$C_{before} = \begin{bmatrix} 1 & 1 & 0.5 & 0 & 0 \\ 1 & 1 & 0.5 & 0 & 0 \\ 0.5 & 0.5 & 1 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 1 & 1 \\ 0 & 0 & 0.5 & 1 & 1 \end{bmatrix} \quad ()$$

%

$$C_{after} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 1 & 1 \\ 0 & 0 & 0.5 & 1 & 1 \end{bmatrix} \quad ()$$

{ }

()

EEAC

$j \quad i$

(ItoU)

$$C(i, j) = \frac{\cap(n_i, n_j)}{\cup(n_i, n_j)} = \frac{n_{i,j}}{n_i + n_j - n_{i,j}} \quad ()$$



•
•

)

(

Iris Wine

¹ UCI repository

² Artificial

() ()

()

	Class	Features	Samples
Glass	6	9	214
Breast-Cancer	2	9	683
Wine	3	13	178
Bupa	2	6	345
Yeast	10	8	1484
Iris	3	4	150
SAHeart	2	9	462
Ionosphere	2	34	351
Halfrings	2	2	400
Galaxy	7	4	323

() Wine

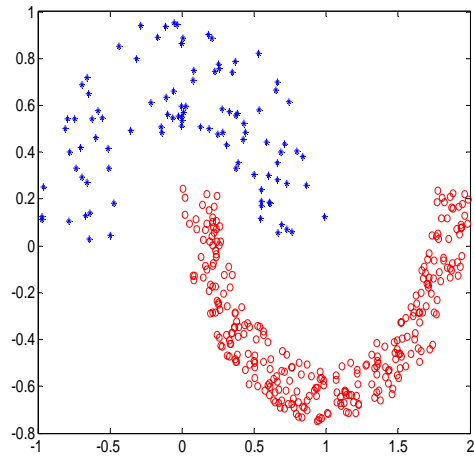
) ()

() ()

Iris

K-means

¹ Newman
² Aeberhard



() " " ()

$N(0,1)$

K-means

K-means

k

k+3 k+2 k+1 k

%

() %

()

()

¹ Relabing

² Single Linkage

N¹.

()

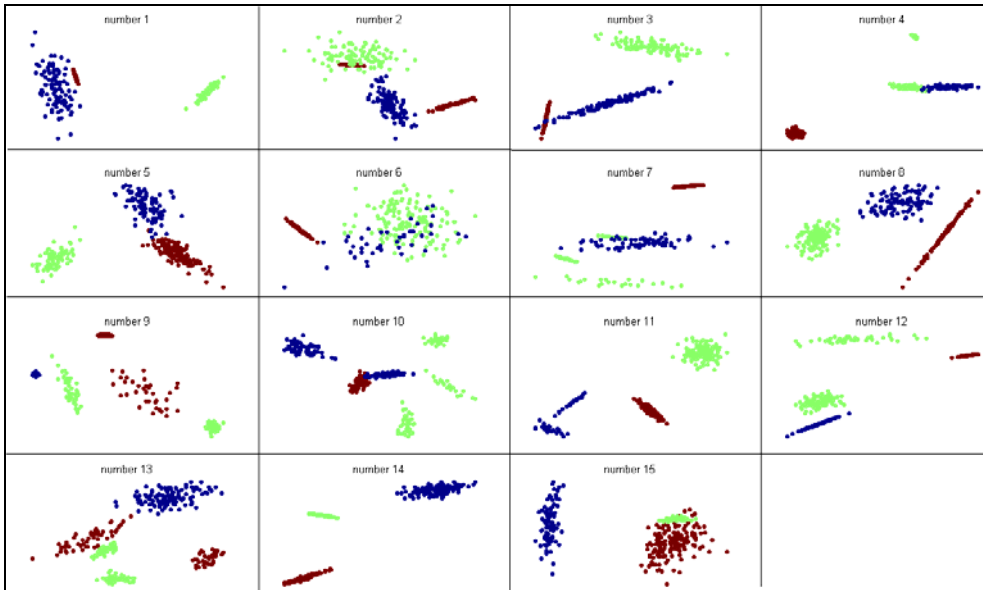
		N. Breast Cancer	Iris	N. Bupa	N. SAHeart	Ionosphere	N. Glass	Halfri ngs	N. Galaxy	N. Yeast	Wine	N. Wine
NMI	ItoU	95.02	88.67	54.78	63.42	70.09	44.86	74.50	29.41	42.86	70.22	96.63
	EEAC	95.73	76.13	54.33	63.36	70.60	47.76	74.48	31.27	42.93	69.38	85.17
MAX	ItoU	96.93	90.00	54.78	64.50	71.51	44.86	87.25	29.41	48.45	71.35	97.75
	EEAC	96.49	84.87	57.42	63.87	57.75	44.35	74.55	29.85	51.27	70.00	94.44
AMM	ItoU	95.43	88.00	54.73	63.42	71.51	44.39	74.50	29.69	48.52	70.73	96.63
	EEAC	95.46	90.00	55.07	63.85	70.66	45.79	54.00	30.65	53.10	70.23	96.63
ENMI	ItoU	96.78	90.00	55.07	64.50	71.51	45.79	88.25	30.03	50.47	70.23	98.32
	EEAC	96.93	88.67	54.78	63.20	71.23	43.93	88.00	30.65	50.47	70.23	97.19
D&Q	1	97.66	97.33	55.36	68.83	72.93	50.47	87.25	35.29	56.67	72.47	98.31
	2	97.07	91.33	55.36	68.02	74.36	53.74	76.50	33.44	54.33	71.35	98.31
	3	97.05	90.00	55.00	63.83	70.91	47.20	83.25	31.36	50.18	71.42	97.75
Adaptive		95.43	88.00	54.73	63.42	71.51	44.39	74.50	29.69	48.52	70.73	96.63
EAC (Full Ens.)		95.17	89.33	54.49	63.20	70.66	46.26	74.50	30.96	44.21	70.22	96.63
Azimi		96.91	89.33	54.75	56.06	70.74	45.05	67.70	29.97	43.40	60.95	96.63

()

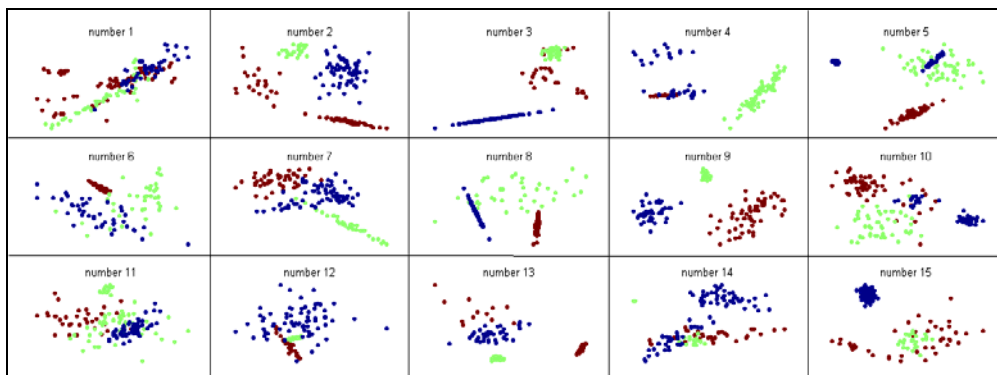
k

¹ Normal

()

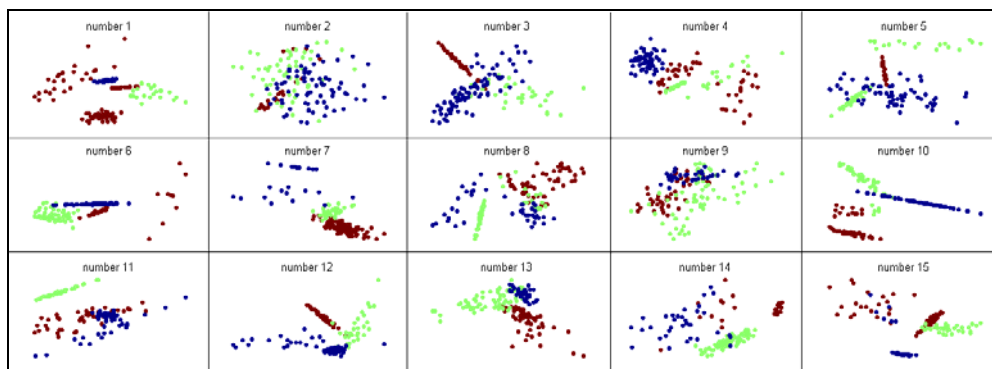


()



()

نتایج و تفسیر آنها



()

»

«

()

1	0.56	94.70	97.67	96.47
2	0.61	61.53	81.27	58.67
3	0.59	88.83	95.27	92.67
4	0.60	69.13	73.00	85.27
5	0.56	88.53	98.00	98.33
6	0.68	81.60	87.33	87.33
7	0.53	77.90	68.57	76.00
8	0.58	90.40	83.03	95.57
9	0.63	72.67	42.33	50.23
10	0.62	60.30	54.67	68.33
11	0.59	81.33	84.20	63.00
12	0.58	76.07	73.73	84.70
13	0.59	77.43	74.43	72.00
14	0.58	84.67	84.67	84.00
15	0.53	89.33	88.90	85.30
Average	0.59	79.63	79.14	79.86

نتایج و تفسیر آنها

»

«

()

1	0.48	54.13	49.60	49.73
2	0.56	59.80	55.47	65.73
3	0.51	87.73	86.47	79.40
4	0.51	55.07	68.60	64.27
5	0.55	74.73	90.00	90.73
6	0.53	79.67	67.73	79.73
7	0.55	57.13	91.93	81.27
8	0.54	72.07	89.20	89.33
9	0.57	100.00	93.87	100.00
10	0.54	66.33	75.20	92.67
11	0.40	41.13	63.73	63.93
12	0.44	72.47	63.20	68.80
13	0.57	92.00	92.00	92.00
14	0.53	59.53	64.00	64.33
15	0.53	66.67	63.33	80.67
Average	0.52	69.23	74.29	77.51

»

«

()

1	0.36	50.93	60.47	55.40
2	0.51	54.73	62.80	56.40
3	0.50	70.87	77.67	74.13
4	0.49	56.07	64.07	66.07
5	0.49	46.67	77.07	77.27
6	0.44	72.47	63.20	68.80
7	0.40	41.13	63.73	63.93
8	0.47	57.53	54.33	47.93
9	0.48	54.13	49.60	49.73
10	0.55	58.53	72.00	72.00
11	0.48	68.93	63.73	82.80
12	0.51	46.33	75.87	68.00
13	0.52	68.60	83.07	88.00
14	0.52	84.73	89.47	89.67
15	0.36	50.93	60.47	55.40
Average	0.48	59.40	68.36	68.58

AMM

D&Q¹

¹ Diversity & Quality

()

%) % (

AMM

()

AMM

%

(ENMI)

%, % ,

(D&Q)

های این روش نسبت به دو روش ترکیب کامل و

روش عظیمی نیز به ترتیب برابر است با ۰.۲ و ۰.۴۲.

%

Aeberhard, S., Coomans D. and de Vel, O. (1992). “*Comparison of Classifiers in High Dimensional Settings*”, *Tech. Rep. no. 92-02, Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland.*

Alizadeh H., Amirgholipour S.K., Seyedaghaee N.R. and Minaei-Bidgoli B. (2009a), *Nearest Cluster Ensemble (NCE): Clustering Ensemble Based Approach for Improving the performance of K-Nearest Neighbor Algorithm*, *11th Conf. of the International Federation of Classification Societies, IFCS09, March 13–18. (in press).*

Alizadeh H., Minaei-Bidgoli B. and Amirgholipour S.K. (2009b), *A New Method for Improving the Performance of K Nearest Neighbor using Clustering Technique*, *International Journal of Convergence Information Technology, JCIT, ISSN: 1975-9320 (in press).*

Ayad H. and Kamel M. (2005), *Cluster-based cumulative ensembles*. In N. Oza and R. Polikar, editors, *Proc. the 6th Intl. Workshop on Multiple Classifier Systems*, pages 236–245. LNCS 3541.

Ayad H.G. and Kamel M.S. (2008), *Cumulative Voting Consensus Method for Partitions with a Variable Number of Clusters*, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, VOL. 30, NO. 1, 160-173.

Barthelemy J.P. and Leclerc B. (1995), *The median procedure for partition*, In *Partitioning Data Sets, AMS DIMACS Series in Discrete Mathematics*, Cox, I. J. et al eds., 19, pp. 3-34.

Baumgartner R., Somorjai R., Summers R., Richter W., Ryner L., and Jarmasz M. (2000), *Resampling as a Cluster Validation Technique in fMRI*, *JOURNAL OF MAGNETIC RESONANCE IMAGING* 11: pp. 228–231.

Banfield C.F. (1976), “*Ultrametric Distances for a Single Linkage Dendrogram*”, *JSTOR: Applied Statistics, Statistical Algorithms*, Vol. 25, No. 3, pp. 313-315.

Ben-Hur A., Elisseeff A. and Guyon I. (2002), *A stability based method for discovering structure in clustered data*, in *Pacific Symposium on Biocomputing*, vol. 7, pp. 6-17.

Brandsma T. and Buishand T.A. (1998), “*Simulation of extreme precipitation in the Rhine basin by nearest-neighbour resampling*”, *Hydrology and Earth System Sciences* 2, pp. 195-209.

Breckenridge J. (1989), *Replicating cluster analysis: Method, consistency and validity*, *Multivariate Behavioral research*.

Brossier G. (1990). *Piecewise hierarchical clustering*, *Journal of Classification*, Springer New York, Vol. 7, No. 2, pp. 197-216.

Das A.K. and Sil J. (2007), *Cluster Validation using Splitting and Merging Technique*, in *proc. of Int. Conf. on Computational Intelligence and Multimedia Applications, ICCIMA*.

Duda R.O., Hart P.E., and Stork D.G. (2001), *Pattern Classification*, second ed. Wiley.

Dudoit S. and Fridlyand, J. (2003), “*Bagging to improve the accuracy of a clustering procedure*”, *Bioinformatics*, 19 (9), pp. 1090-1099.

Estivill-Castro V. and Yang J. (2003), *Cluster Validity Using Support Vector Machines*, *DaWaK 2003, LNCS 2737*, pp. 244–256.

Faceli K., Marcilio C.P. Souto d. (2006), *Multi-objective Clustering Ensemble*, *Proceedings of the Sixth International Conference on Hybrid Intelligent Systems (HIS'06)*.

-
- Fern, X. and Brodley, C. E. (2003). *Random Projection for High Dimensional Data Clustering: A Cluster Ensemble Approach*, In Proc. 20th Int. conf. on Machine Learning, ICML 2003.
- Fern X. and Lin W. (2008), *Cluster Ensemble Selection*, SIAM International Conference on Data Mining (SDM08).
- Fischer B. and Buhmann J.M., (2003), *Bagging for path-based clustering*, IEEE Transactions on Pattern Analysis and Machine Intelligence, pp.1411–1415.
- Fred, A. and Jain, A. K. (2002). “*Data Clustering Using Evidence Accumulation*”, Proc. of the 16th Intl. Conf. on Pattern Recognition, ICPR02, Quebec City, pp. 276 – 280.
- Fred A. and Jain A.K., (2003), *Robust data clustering*, in: Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR, USA, vol. II, pp. 128–136.
- Fred A.L. and Jain A.K. (2005). *Combining Multiple Clusterings Using Evidence Accumulation*. IEEE Trans. on Pattern Analysis and Machine Intelligence, 27(6):835–850.
- Fred A. and Jain A.K. (2006), *Learning Pairwise Similarity for Data Clustering*, In Proc. of the 18th Int. Conf. on Pattern Recognition (ICPR’06).
- Fred A. and Lourenco A. (2008), *Cluster Ensemble Methods: from Single Clusterings to Combined Solutions*, Studies in Computational Intelligence (SCI), 126, 3–30.
- Fridlyand J. and Dudoit S. (2001). *Applications of resampling methods to estimate the number of clusters and to improve the accuracy of a clustering method*. Stat. Berkeley Tech Report. No. 600.
- Jain A., Murty M. N., and Flynn P. (1999), *Data clustering: A review*. ACM Computing Surveys, 31(3):264–323.
- Inokuchi R., Nakamura T. and Miyamoto S. (2006), *Kernelized Cluster Validity Measures and Application to Evaluation of Different Clustering Algorithms*, in proc. of the IEEE Int. Conf. on Fuzzy Systems, Canada, July 16-21.
- Jain A.K. and Dubes R.C. (1988), *Algorithms for Clustering Data*. Prentice Hall.
- Kaufman L. and Rosseeuw P.J. (1990), *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley & Sons, Inc.
- Kuncheva L.I. and Hadjitodorov S. (2004). *Using diversity in cluster ensembles*. In Proc. of IEEE Intl. Conference on Systems, Man and Cybernetics, pages 1214–1219.
- Kuncheva L.I. and Whitaker C. J., (2003), *Measures of diversity in classifier ensembles*, Machine Learning.
- Lange T., Braun M.L., Roth V., and Buhmann J.M. (2003). *Stability-based model selection*. In Advances in Neural Information Processing Systems 15. MIT Press.
- Lange T., Roth V., Braun M.L., and Buhmann J.M. (2004). *Stability-based validation of clustering solutions*. Neural Computation, 16(6):1299–1323.
- Lapointe F.J. and Legendre P. (1991). *The generation of random ultrametric matrices representing dendrograms*. Journal of Classification, Springer New York, Vol. 8, No. 2, pp 177–200.
- Law M.H.C., Topchy A.P., and Jain A.K. (2004). *Multiobjective data clustering*. In Proc. of IEEE Conference on Computer Vision and Pattern Recognition, volume 2, pages 424–430, Washington D.C.
- Levine E., Domany E. (2001), *Resampling Method for Unsupervised Estimation of Cluster Validity*. Neural Computation 13: 2573-2593.

- Luxburg U.V. and Ben-David S. (2005), *Towards a statistical theory of clustering*, Technical report, PASCAL workshop on clustering, London.
- Man Y. and Gath I. (1994), “*Detection and Separation of Ring-Shaped Clusters Using Fuzzy Clusters*” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 16, no. 8, pp. 855-861.
- Minaei-Bidgoli B., Topchy A. and Punch W.F. (2004), “*Ensembles of Partitions via Data Resampling*”, in *Proc. Intl. Conf. on Information Technology, ITCC 04, Las Vegas*.
- Mohammadi M., Alizadeh H. and Minaei-Bidgoli B. (2008), “*Neural Network Ensembles using Clustering Ensemble and Genetic Algorithm*”, *Intl. Conf. on Convergence and hybrid Information Technology, ICCIT08, Nov. 11-13, IEEE CS*.
- Moller U., Radke D. (2006), *A Cluster Validity Approach based on Nearest-Neighbor Resampling*, In *Proc. of the 18th Int. Conf. on Pattern Recognition (ICPR'06)*.
- Newman C.B.D.J., Hettich S. and Merz C. (1998), *UCI repository of machine learning databases*, <http://www.ics.uci.edu/~mllearn/MLSummary.html>.
- Parvin H., Alizadeh H. and Minaei-Bidgoli B. (2009a), *A New Method for Constructing Classifier Ensembles*, *International Journal of Digital Content: Technology and its Application, JDCTA, ISSN: 1975-9339, (in press)*.
- Parvin H., Alizadeh H. and Minaei-Bidgoli B. (2009b), *Using Clustering for Generating Diversity in Classifier Ensemble*, *International Journal of Digital Content: Technology and its Application, JDCTA, ISSN: 1975-9339, Vol. 3, No.1, pp. 51-57*.
- Roth V., Lange T., Braun M., and Buhmann J. (2002), *A Resampling Approach to Cluster Validation*, *Intl. Conf. on Computational Statistics, COMPSTAT*.
- Rakhlin A. and Caponnetto A. (2007), *Stability of k-means clustering*, In *Advances in Neural Information Processing Systems 19, MIT Press, Cambridge, MA*.
- Roth V. and Lange T. (2004), *Feature Selection in Clustering Problems*, In *Advances in Neural Information Processing Systems, NIPS04*.
- Roth V., Braun M.L., Lange T., and Buhmann J.M. (2002), *Stability-Based Model Order Selection in Clustering with Applications to Gene Expression Data*, *ICANN 2002, LNCS 2415, pp. 607–612*.
- Shamiry O., Tishby N. (2007), *Cluster Stability for Finite Samples*, *21st Annual Conference on Neural Information Processing Systems (NIPS07)*.
- Strehl A. and Ghosh J. (2002), *Cluster ensembles - a knowledge reuse framework for combining multiple partitions*. *Journal of Machine Learning Research*, 3(Dec):583–617.
- Topchy, A., Jain, A.K. and Punch, W.F. (2003), “*Combining Multiple Weak Clusterings*”, *Proc. 3d IEEE Intl. Conf. on Data Mining*, pp. 331-338.
- Topchy A., Minaei-Bidgoli B., Jain A.K. and Punch W.F., (2004), *Adaptive Clustering ensembles*, In *Proc. Intl. Conf on Pattern Recognition, ICPR'04, Cambridge, UK, pp.272-275*.
- Xie X.L., Beni G. (1991), *A Validity measure for Fuzzy Clustering*, *IEEE Trans. on Pattern Analysis and Machine Intelligence*, Vol.13, No.4, pp. 841–846.

" *K-means*

" ()

()

() .

() .

() .

() .

Abstract:

Data clustering is an important step in data mining which excavates hidden patterns in unsupervised data. On account of complexity of the problem and weakness of the base clustering algorithms, most of researches are focused on ensemble methods. Diversity in primary results and also quality of the primary results are two vital factors which affect the final results. Although both of these two factors are considered and investigated in the late researches on clustering ensembles; there are some of questions which are still vague. Using a subset of primary results can be better than total primary results or not? Which subset of the primary results can cause to improve the performance of clustering ensemble? How to evaluate the primary results? This thesis tries to give a reasonable answer to these questions. Here, a framework for improving the efficiency of the clustering ensemble is proposed which is based on using a subset of primary clusters. Furthermore, some new methods for each step of this framework are suggested. For evaluating each individual cluster some new methods are proposed which are inspired from mutual information. Moreover, two new techniques are presented to construct the co-association matrix from only a subset of primary clusters. Experimental results over several standard data sets show that the presented schemes can significantly enhance the efficiency of the primary and even the full ensemble results. The average enhancement over 11 examined cases is 2.3% in comparison with the full ensemble. In addition, several artificial data sets are produced and examined during empirical studies.

Keywords:

Clustering Ensemble, Cluster Evaluation, Mutual Information, Subset of Primary Results, Evidence Accumulation Clustering, Co-association Matrix.



**Iran University of Science and Technology
Computer Engineering Department**

Clustering Ensemble Based on a Subset of Primary Results

**A Thesis Submitted in Partial Fulfillment of the Requirement for the
Degree of Master of Science in Computer Engineering**

**By:
Hosein Alizadeh**

**Supervisor:
Dr. Behrouz Minaei-Bidgoli**

March 2009