

Improving Cluster Analysis by Co-initializations (Supplemental Document)

He Zhang*, Zhirong Yang, Erkki Oja

Department of Information and Computer Science, Aalto University, Espoo, Finland

Abstract

In this document we give the following supplemental information: 1) detailed description of the datasets in the experiments, 2) the learning objective and constraints of the compared methods, 3) more experiment results that are complementary to those in the paper.

1. Datasets

Table 1 gives the information (sources and links) of 19 datasets used in the experiments. These datasets have been widely utilized by machine learning and data mining community, and are freely downloadable from the Internet. A brief description is given in the following (problems, samples, classes, and dimensionalities).

- **ORL**: the AT&T *ORL* database of face images. There are 10 different images for each of 40 distinct subjects with varying lighting conditions and facial expressions. The size of each image is 92×112 pixels, with 256 grey levels per pixel. In summary, the dataset contains 400 samples grouped in 40 classes, with each sample having a dimensionality of 10,304.
- **MED**: the *MED* database contains abstract text collections. There are 696 documents organized in 25 topics, with a dictionary containing 5,831 words. In summary, the dataset contains 696 samples grouped in 25 classes, with each sample having a dimensionality of 5,831.
- **VOWEL**: the LIBSVM *vowel* dataset, originally from UCI machine learning repository. The problem is specified by the accompanying data file, “vowel.data” which consists of a three dimensional array: `voweldata [speaker, vowel, input]`. The speakers are indexed by integers 0-89. The vowels are indexed by integers 0-10. For each utterance, there are ten floating-point input values, with array

*Corresponding author.

Email addresses: `he.zhang@aalto.fi` (He Zhang), `zhirong.yang@aalto.fi` (Zhirong Yang), `erkki.oja@aalto.fi` (Erkki Oja)

Table 1: Data information.

Dataset	Source	URL
ORL	ORL	http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html
MED	LSI	http://web.eecs.utk.edu/research/lsi/
VOWEL	UCI	http://archive.ics.uci.edu/ml/
COIL20	COIL	http://www.cs.columbia.edu/CAVE/software/softlib/coil-20.php
SEMEION	UCI	http://archive.ics.uci.edu/ml/
FAULTS	UCI	http://archive.ics.uci.edu/ml/
SEGMENT	UCI	http://archive.ics.uci.edu/ml/
CORA	LINQS	http://www.cs.umd.edu/projects/linqs/projects/lbc/index.html
CITeseer	LINQS	http://www.cs.umd.edu/projects/linqs/projects/lbc/index.html
7SECTORS	CMUTE	http://www.cs.cmu.edu/~TextLearning/datasets.html
OPTDIGITS	UCI	http://archive.ics.uci.edu/ml/
SVMGUIDE1	LIBSVM	http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/
ZIP	LIBSVM	http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/
USPS	UCI	http://archive.ics.uci.edu/ml/
PENDIGITS	UCI	http://archive.ics.uci.edu/ml/
PROTEIN	LIBSVM	http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/
20NEWS	CMUTE	http://www.cs.cmu.edu/~TextLearning/datasets.html
LET-REC	UCI	http://archive.ics.uci.edu/ml/
MNIST	MNIST	http://yann.lecun.com/exdb/mnist/

20 indices 0-9. The problem is to train the network as well as possible using only
 21 on data from “speakers” 0-47, and then to test the network on speakers 48-89,
 22 reporting the number of correct classifications in the test set. In summary, the
 23 dataset contains 990 samples grouped in 11 classes, with each sample having a
 24 dimensionality of 10.

25 • COIL20: the *COIL-20* dataset is a toy image database from Columbia University
 26 Image Library. It contains 1,440 grayscale images of 20 objects (72 images per
 27 object) with a wide variety of complex geometric and reflectance characteristics.
 28 Each image has a size of 128×128 . The dataset has been widely used in image
 29 classification and retrieval tasks. In summary, the dataset contains 1,440 samples
 30 grouped in 20 classes, with each sample having a dimensionality of 16,384.

31 • SEMEION: the UCI *Semeion Handwritten Digit* dataset. Totally 1,593 hand-
 32 written digits from around 80 persons were scanned, stretched in a rectangular
 33 box 16×16 in a gray scale of 256 values. Then each pixel of each image was
 34 scaled into a boolean (1/0) value using a fixed threshold. In summary, the dataset
 35 contains 1,593 samples grouped in 10 classes, with each sample having a dimen-

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73

sionality of 256.

- FAULTS: the UCI *Steel Plates Faults* dataset consists of information of steel plate faults, classified into 7 different types, with 27 independent attributes. The goal was to train machine learning algorithms for automatic pattern recognition. In summary, the dataset contains 1,941 samples grouped in 7 classes, with each sample having a dimensionality of 27.
- SEGMENT: the UCI *Image Segmentation* dataset. The instances, represented by 19 high-level features, were drawn randomly from a database of 7 outdoor images, and the images were handsegmented to create a classification for every pixel. In summary, the dataset contains 2,310 samples grouped in 7 classes, with each sample having a dimensionality of 19.
- CORA: the LINQS *Cora* dataset. It consists of 2,708 scientific publications classified into one of 7 classes. The citation network consists of 5,429 links. Each publication in the dataset is described by a 0/1-valued word vector indicating the absence/presence of the corresponding word from the dictionary. The dictionary consists of 1,433 unique words. In summary, the dataset contains 2,708 samples grouped in 7 classes, with each sample having a dimensionality of 1,433.
- CITeseer: the LINQS *CiteSeer* dataset. It consists of 3,312 scientific publications classified into one of 6 classes. The citation network consists of 4,732 links. Each publication in the dataset is described by a 0/1-valued word vector indicating the absence/presence of the corresponding word from the dictionary. The dictionary consists of 3,703 unique words. In summary, the dataset contains 3,312 samples grouped in 6 classes, with each sample having a dimensionality of 3,703.
- 7SECTORS: the *Industry Sector* dataset from CMU Text Learning group. The dataset is a collection of web pages belonging to companies from 7 economic sectors. In summary, the dataset contains 4,556 samples grouped in 7 classes, with each sample having a dimensionality of 10,000.
- OPTDIGITS: the UCI *optical recognition of handwritten digits* dataset was created by extracting normalized bitmaps of handwritten digits of 43 people from a preprinted form, which generates an input matrix of 8×8 with each element being an integer in the range $[0, 16]$. In summary, the dataset contains 5,620 samples grouped 10 classes, with each sample having a dimensionality of 64.
- SVMGUIDE1: the LIBSVM *svmguide1* dataset is obtained from an astroparticle application from Jan Conrad of Uppsala University, Sweden. There are 3,089 instances for training and 4,000 for testing, with each instance represented by 4 numerical features. In the paper, we used all the samples, i.e., 7,089 samples grouped in 2 classes, with each sample having a dimensionality of 4.

- 74 • ZIP: the LIBSVM *ZIP* handwritten digits dataset contains 9,298 samples grouped
75 10 classes, with each sample having a dimensionality of 256.
 - 76 • USPS: the UCI *optical recognition of handwritten digits* dataset is used for op-
77 tical character recognition, similar to OPTDIGITS. It contains 9,298 samples
78 grouped 10 classes, with each sample having a dimensionality of 256.
 - 79 • PENDIGITS: the UCI *pen-based recognition of handwritten digits* dataset was
80 created as a digit database that contains 250 samples from 44 writers. It contains
81 10,992 samples grouped in 10 classes, with each sample having a dimensionality
82 of 16.
 - 83 • PROTEIN: the LIBSVM *protein* dataset from bioinformatics. The original dataset
84 has 17,766 instances for training and 6,621 for testing, with each sample repre-
85 sented by 357 features. In the paper, we utilized the training subset, i.e., 17,766
86 samples grouped in 3 classes, with each sample having a dimensionality of 357.
 - 87 • 20NEWS: text documents from *20 newsgroups*. This data set is a collection of
88 20,000 messages, collected from 20 different netnews newsgroups. One thou-
89 sand messages from each of the twenty newsgroups were chosen at random and
90 partitioned by newsgroup name. 10,000 words with maximum information gain
91 are preserved. The dataset we used in the paper contains 19,938 samples grouped
92 in 20 classes, with each sample having a dimensionality of 10,000.
 - 93 • LET-REC: the UCI *letter recognition* dataset. The objective is to identify each
94 of a large number of black-and-white rectangular pixel displays as one of the 26
95 capital letters in the English alphabet. The character images were based on 20
96 different fonts and each letter within these 20 fonts was randomly distorted to
97 produce a file of 20,000 unique stimuli. Each stimulus was converted into 16
98 primitive numerical attributes (statistical moments and edge counts) which were
99 then scaled to fit into a range of integer values from 0 through 15. In summary,
100 the dataset contains 20,000 samples grouped 26 classes, with each sample having
101 a dimensionality of 16.
 - 102 • MNIST: the handwritten digit images database. The MNIST database has a train-
103 ing set of 60,000 examples, and a test set of 10,000 examples. It is a subset of
104 a larger set available from NIST. The digits have been size-normalized and cen-
105 tered in a 28×28 image. In the paper, we used all the samples, i.e., 70,000
106 samples grouped in 10 classes, with each sample having a dimensionality of 784.
- 107 Note that, as a pre-processing step, the scattering features [1] have been extracted for
108 each sample in the image datasets, and Tf-Idf features have been extracted for the text
109 document datasets.

110 **2. Optimization specifications**

111 In this section we give the optimization specifications of the compared methods.
 112 Note that the specification of DCD has been presented in the paper already. Table 2 gives the frequently used notations.

Table 2: List of frequently used notations.

$G = (V, E)$	undirected graph G with vertices in V and edges in E
m, n, r	data dimensionality, sample size, reduced rank of matrix
$\mathbb{R}_+^{m \times n}$	space of nonnegative $m \times n$ matrices
X	data matrix of size $m \times n$, whose columns are n -dimensional vectors
A	similarity matrix of size $n \times n$
W	factorizing matrix of size $n \times r$, also called cluster indicator matrix
S	factorizing matrix of size $r \times r$
$\Pi = \{\pi_1, \dots, \pi_M\}$	a cluster ensemble with M base clusterings (partitions)

113

- **KM:** the classical K-means algorithm [2]. Let $X = \{x_i, i = 1, \dots, n$ be the set of m -dimensional points to be clustered into a set of K clusters, $C = \{c_k, k = 1, \dots, K\}$. K-means algorithm finds a partition that minimizes the squared error between the empirical mean of a cluster and the points belonging to the cluster. The aim of K-means is to minimize the sum of the squared error over all clusters, given by the objective function:

$$J(C) = \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - \mu_k\|^2,$$

114 where μ_k is the empirical mean of the cluster C_k . We directly utilized the
 115 Matlab function *kmeans* for implementation.

- **NCUT:** Normalized Cut [3]. NCUT partitions the graph G into two disjoint sets A and B by minimizing the cost as a fraction of the total edge connections to all the nodes in the graph, given by the objective function:

$$NCUT(A, B) = \frac{cut(A, B)}{assoc(A, V)} + \frac{cut(A, B)}{assoc(B, V)},$$

116 where $cut(A, B) = \sum_{u \in A, v \in B} w(u, v)$ denotes the degree of dissimilarity between
 117 the subgraphs A and B computed as total weight of the edges removed for sep-
 118 arating the two subgraphs, and $assoc(A, V) = \sum_{u \in A, t \in V} w(u, t)$ is the total con-
 119 nection from nodes in A to all nodes in the graph and $assoc(B, V)$ is similarly
 120 defined. The objective given above can be minimized by solving a generalized
 121 eigenvalue problem (see details in [3]). For implementation, we utilized the

122

NCUT Matlab package¹ downloaded from the author's website.

- 1-SPEC: 1-Spectral Ratio Cheeger Cut [4]. The ratio Cheeger cut (RCC) of a partition (C, \bar{C}) , where $C \subset V$ and $\bar{C} = V \setminus C$, is to minimize the objective function:

$$RCC(C, \bar{C}) = \frac{cut(C, \bar{C})}{\min\{|C|, |\bar{C}|\}}.$$

123

We adopted the 1-SPEC software² by Hein and Bühler with its default setting in our implementation.

124

125

126

127

128

129

- PNMF: Projective NMF [5]. Given the nonnegative input data matrix $X \in \mathbb{R}_+^{n \times m}$, PNMF [6, 7] based on Frobenius norm aims to find a factorizing matrix $W \in \mathbb{R}_+^{n \times r}$ for the optimization problem: minimize $\|X - WW^T X\|_F^2$. In our paper, we utilized the kernelized version of PNMF, i.e., to replace XX^T by the similarity matrix A between samples.

- NSC: Nonnegative Spectral Clustering [8]. NSC solves the normalized cut by using the multiplicative update algorithm of NMF, i.e., to solve the following optimization problem:

$$\underset{H^T D H = I, H \geq 0}{\text{minimize}} \quad -\text{trace}(H^T A H),$$

130

where $D = \text{diag}(d_1, \dots, d_n)$ with $d_i = \sum_{j=1}^n A_{i,j}$.

- ONMF: Symmetric Tri-Factor Orthogonal NMF [9]. ONMF is the special case of orthogonal tri-factor NMF, when the given input is a matrix of pairwise similarities. ONMF thus solves the following optimization problem:

$$\underset{W \geq 0, S \geq 0}{\text{minimize}} \quad \|A - WS W^T\|_F^2, \quad \text{subject to } W^T W = I,$$

131

where $S \in \mathbb{R}_+^{r \times r}$.

- LSD: Left Stochastic Matrix Decomposition [10]. LSD is a probabilistic clustering method. It estimates a scaling factor c^* and a cluster probability matrix W^*

¹<http://www.cis.upenn.edu/~jshi/software/>

²<http://www.ml.uni-saarland.de/code/oneSpectralClustering/oneSpectralClustering.html>

that solves the following optimization problem:

$$\underset{c \in \mathbb{R}_+}{\text{minimize}} \left\{ \underset{W \geq 0}{\text{minimize}} \|cA - WW^T\|_F^2, \text{ subject to } \sum_{k=1}^r W_{ik} = 1 \right\}.$$

132 Note that minimizing the scaling factor c^* is given in a closed form and does
 133 not depend on a particular solution W^* , which means that only W needs to be
 134 updated.

- PLSI: Probabilistic Latent Semantic Indexing [11]. PLSI assumes that the data is generated from a multinomial distribution, and maximizing the PLSI likelihood function is equivalent to minimizing the Kullback-Leibler divergence. We utilized the symmetric version of PLSI in our context, i.e., to solve the following optimization problem:

$$\underset{W \geq 0, S \geq 0}{\text{minimize}} \left\{ A_{ij} \log \frac{A_{ij}}{(WSW^T)_{ij}} \right\}, \text{ subject to } \sum_{i=1}^n W_{ik} = 1 \text{ and } \sum_{k=1}^r S_{kk} = 1,$$

135 where S is a diagonal matrix.

- BEST: BESTCLUSTERING ensemble algorithm [12]. Given a set of n samples $X = \{x_i\}$ and a cluster ensemble of M base clusterings $\Pi = \{\pi_1, \dots, \pi_M\}$, the following simple 0/1 distance function checks if two clusterings π_1 and π_2 place x_i and x_j in the same clusters:

$$d_{(x_i, x_j)}(\pi_1, \pi_2) = \begin{cases} 1 & \text{if } \pi_1(x_i) = \pi_2(x_j) \text{ and } \pi_1(x_i) \neq \pi_2(x_j), \\ & \text{or } \pi_1(x_i) \neq \pi_2(x_j) \text{ and } \pi_1(x_i) = \pi_2(x_j), \\ 0 & \text{otherwise.} \end{cases}$$

The distance between two clusterings π_1 and π_2 is defined as the number of pairs of objects on which the two clusterings disagree, that is,

$$d_X(\pi_1, \pi_2) = \sum_{(x_i, x_j) \in X \times X} d_{(x_i, x_j)}(\pi_1, \pi_2)$$

Therefore, BEST algorithm finds the clustering π^* from Π that minimizes the total number of disagreements with all the given clusterings, i.e., to solve the following optimization problem:

$$\underset{\pi^*}{\text{minimize}} \left\{ \sum_{i=1}^M d_X(\pi_i, \pi^*), \pi^* \neq \pi_i \right\}.$$

- CO: evidence accumulation method based on CO-association matrix [13]. The assumption is that patterns belonging to a “natural” cluster are very likely to

be co-located in the same cluster in different data partitions π_i . Taking the co-occurrences of pairs of patterns in the same cluster as votes for their association, the M data partitions of n patterns are mapped into a $n \times n$ co-association matrix:

$$CO(i, j) = \frac{n_{ij}}{M},$$

136 where where n_{ij} is the number of times the pattern pair (x_i, x_j) is assigned to the
 137 same cluster among the M partitions. The evidence accumulation mechanism
 138 thus maps the partitions in the clustering ensemble into a new similarity mea-
 139 sure (the co-association matrix CO) between patterns. Therefore, any similarity-
 140 based clustering algorithms can be used to produce the final clustering result.
 141 In this paper, we obtained the final partition by using the complete-linkage hi-
 142 erarchical clustering algorithm (implemented by the authors in [14] as a Matlab
 143 package³).

- CTS: link-based ensemble clustering with Connected-Triple based Similarity matrix [15]. CTS method can better handle the unknown relations between data samples than the CO method does. Briefly, the weight assigned to edge connecting clusters i and j is estimated in accordance with the proportion of their overlapping members $w_{ij} = \frac{|X_i \cap X_j|}{|X_i \cup X_j|}$, where X_C denotes set of data points belonging to cluster C . The count of all triples $1, \dots, q$ between cluster i and cluster j can be calculated as $C_{ij} = \sum_{k=1}^q \{min(w_{ik}, w_{jk})\}$. The similarity between clusters i and j can be estimated as follows:

$$S_{WT}(i, j) = \frac{C_{ij}}{C_{max}},$$

where C_{max} is the maximum C_{ij} value of any two clusters i and j . As for the CTS method, for the m -th ensemble member, the similarity of data sample x_i and x_j is estimated as:

$$S_m(x_i, x_j) = \begin{cases} 1 & \text{if } C(x_i) = C(x_j), \\ S_{WT}(C(x_i), C(x_j)) \times DC & \text{otherwise.} \end{cases}$$

where DC is a constant decay factor ranging in $[0, 1]$ (i.e. the confidence level of accepting two non-identical objects as being similar). Finally, each entry in the CTS similarity matrix can be computed as:

$$CTS(x_i, x_j) = \frac{1}{M} \sum_{m=1}^M S_m(x_i, x_j).$$

144 Note that we directly ran the CTS algorithm from the Matlab package mentioned

³<http://www.jstatsoft.org/v36/i09>

145 in [14], and produced the final partition using the complete-linkage algorithm as
146 well.

147 3. Experiments

148 Table 3 gives the complete clustering results of the first group of experiments (see
149 Section 5.3 of the original paper). Table 4 shows the clustering results of the second
150 group of experiments (see Section 5.3 of the original paper), where the ensemble bases
151 are created by the classical k -means algorithm.

152 References

- 153 [1] S. Mallat, Group invariant scattering, *Communications on Pure and Applied*
154 *Mathematics* 65 (10) (2012) 1331–1398.
- 155 [2] S. Lloyd, Last square quantization in pcm, *IEEE Transactions on Information*
156 *Theory*, special issue on quantization 28 (1982) 129–137.
- 157 [3] J. Shi, J. Malik, Normalized cuts and image segmentation, *IEEE Transactions on*
158 *Pattern Analysis and Machine Intelligence* 22 (8) (2000) 888–905.
- 159 [4] M. Hein, T. Bühler, An inverse power method for nonlinear eigenproblems with
160 applications in 1-Spectral clustering and sparse PCA, in: *Advances in Neural*
161 *Information Processing Systems (NIPS)*, 2010, pp. 847–855.
- 162 [5] Z. Yang, E. Oja, Linear and nonlinear projective nonnegative matrix factorization,
163 *IEEE Transactions on Neural Networks* 21 (5) (2010) 734–749.
- 164 [6] Z. Yuan, E. Oja, Projective nonnegative matrix factorization for image compres-
165 sion and feature extraction, in: *Proceedings of 14th Scandinavian Conference on*
166 *Image Analysis (SCIA)*, Joensuu, Finland, 2005, pp. 333–342.
- 167 [7] Z. Yang, Z. Yuan, J. Laaksonen, Projective non-negative matrix factorization with
168 applications to facial image processing, *International Journal on Pattern Recog-
169 nition and Artificial Intelligence* 21 (8) (2007) 1353–1362.
- 170 [8] C. Ding, T. Li, M. Jordan, Nonnegative matrix factorization for combinatorial
171 optimization: Spectral clustering, graph matching, and clique finding, in: *IEEE*
172 *International Conference on Data Mining (ICDM)*, IEEE, 2008, pp. 183–192.
- 173 [9] C. Ding, T. Li, W. Peng, H. Park, Orthogonal nonnegative matrix t-factorizations
174 for clustering, in: *International Conference on Knowledge Discovery and Data*
175 *Mining (SIGKDD)*, ACM, 2006, pp. 126–135.
- 176 [10] R. Arora, M. Gupta, A. Kapila, M. Fazel, Clustering by left-stochastic matrix
177 factorization, in: *International Conference on Machine Learning (ICML)*, 2011,
178 pp. 761–768.

- 179 [11] T. Hofmann, Probabilistic latent semantic indexing, in: International Conference
180 on Research and Development in Information Retrieval (SIGIR), 1999, pp. 50–
181 57.
- 182 [12] A. Gionis, H. Mannila, P. Tsaparas, Clustering aggregation, in: International Con-
183 ference on Data Engineering (ICDE), IEEE, 2005, pp. 341–352.
- 184 [13] A. Fred, A. Jain, Combining multiple clusterings using evidence accumulation,
185 IEEE Transactions on Pattern Analysis and Machine Intelligence 27 (6) (2005)
186 835–850.
- 187 [14] N. Iam-On, S. Garrett, Linkclue: A matlab package for link-based cluster ensem-
188 bles, Journal of Statistical Software 36 (9) (2010) 1–36.
- 189 [15] N. Iam-On, T. Boongoen, S. Garrett, Refining pairwise similarity matrix for clus-
190 ter ensemble problem with cluster relations, in: International Conference on Dis-
191 covery Science (DS), Springer, 2008, pp. 222–233.

DATASET	KM	NCUT	1-SPEC	PNMF	NSC	ONMF	LSD	PLSI	DCD1	DCD1.2	DCD2	DCD5
ORL	0.70	0.81	0.81	0.81	0.81	0.81	0.81	0.81	0.81	0.81	0.81	0.81
MED	0.59	0.57	0.53	0.57	0.59	0.57	0.57	0.57	0.57	0.57	0.57	0.57
VOWEL	0.40	0.35	0.34	0.38	0.35	0.35	0.35	0.35	0.35	0.35	0.35	0.35
COIL20	0.63	0.71	0.67	0.67	0.71	0.72	0.72	0.71	0.68	0.68	0.68	0.68
SEMION	0.68	0.64	0.66	0.60	0.66	0.60	0.60	0.60	0.60	0.60	0.60	0.60
FAULTS	0.42	0.40	0.40	0.42	0.45	0.45	0.45	0.45	0.45	0.45	0.45	0.45
SEGMENT	0.59	0.61	0.55	0.49	0.54	0.49	0.53	0.53	0.48	0.51	0.51	0.51
CORA	0.53	0.39	0.36	0.41	0.36	0.36	0.36	0.36	0.41	0.41	0.41	0.41
CITFSEER	0.61	0.30	0.31	0.28	0.29	0.28	0.28	0.28	0.30	0.30	0.30	0.30
7SECTORS	0.39	0.25	0.25	0.29	0.29	0.29	0.29	0.29	0.30	0.30	0.30	0.30
OPTDIGITS	0.72	0.74	0.76	0.70	0.68	0.68	0.68	0.68	0.71	0.71	0.71	0.71
SVMGUIDE1	0.71	0.75	0.93	0.68	0.68	0.68	0.68	0.68	0.70	0.70	0.70	0.70
ZIP	0.49	0.74	0.74	0.54	0.70	0.72	0.68	0.68	0.57	0.57	0.57	0.57
USPS	0.74	0.74	0.74	0.67	0.80	0.75	0.68	0.68	0.72	0.74	0.74	0.74
PENNDIGITS	0.72	0.80	0.73	0.79	0.79	0.79	0.79	0.79	0.80	0.80	0.80	0.80
PROTEIN	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46	0.46
20NEWS	0.07	0.43	0.36	0.39	0.39	0.39	0.38	0.38	0.41	0.41	0.41	0.41
LEFT-REC	0.29	0.21	0.15	0.36	0.37	0.34	0.35	0.35	0.17	0.21	0.21	0.21
MINST	0.60	0.77	0.88	0.57	0.87	0.73	0.57	0.57	0.46	0.79	0.97	0.97

(a) Purity

DATASET	KM	NCUT	1-SPEC	PNMF	NSC	ONMF	LSD	PLSI	DCD1	DCD1.2	DCD2	DCD5
ORL	0.85	0.90	0.92	0.89	0.90	0.90	0.90	0.90	0.83	0.90	0.91	0.91
MED	0.55	0.57	0.52	0.56	0.55	0.55	0.55	0.55	0.56	0.57	0.57	0.57
VOWEL	0.43	0.40	0.38	0.39	0.36	0.37	0.39	0.39	0.28	0.39	0.38	0.41
COIL20	0.77	0.79	0.77	0.75	0.79	0.79	0.79	0.79	0.74	0.80	0.80	0.80
SEMION	0.57	0.61	0.62	0.58	0.62	0.58	0.58	0.58	0.61	0.61	0.61	0.61
FAULTS	0.10	0.08	0.09	0.09	0.11	0.11	0.11	0.11	0.02	0.08	0.11	0.11
SEGMENT	0.58	0.55	0.58	0.43	0.48	0.43	0.49	0.49	0.07	0.55	0.53	0.58
CORA	0.34	0.16	0.14	0.14	0.13	0.17	0.17	0.17	0.15	0.20	0.24	0.26
CITFSEER	0.34	0.10	0.12	0.07	0.07	0.07	0.07	0.07	0.10	0.18	0.21	0.21
7SECTORS	0.17	0.04	0.05	0.05	0.04	0.04	0.04	0.04	0.04	0.08	0.13	0.11
OPTDIGITS	0.70	0.72	0.80	0.67	0.68	0.67	0.67	0.67	0.69	0.69	0.69	0.69
SVMGUIDE1	0.31	0.35	0.65	0.27	0.27	0.27	0.27	0.27	0.02	0.40	0.39	0.59
ZIP	0.40	0.78	0.79	0.54	0.67	0.65	0.64	0.64	0.21	0.78	0.79	0.81
USPS	0.62	0.77	0.80	0.60	0.75	0.71	0.66	0.66	0.46	0.76	0.77	0.81
PENNDIGITS	0.68	0.81	0.78	0.78	0.78	0.78	0.78	0.78	0.10	0.81	0.83	0.86
PROTEIN	0.00	0.01	0.01	0.02	0.01	0.01	0.01	0.01	0.02	0.01	0.04	0.04
20NEWS	0.05	0.54	0.52	0.36	0.36	0.34	0.34	0.34	0.14	0.45	0.44	0.45
LEFT-REC	0.35	0.38	0.26	0.43	0.43	0.43	0.42	0.42	0.18	0.36	0.42	0.46
MINST	0.51	0.81	0.89	0.59	0.82	0.72	0.59	0.59	0.48	0.80	0.92	0.93

(b) NMI

DATASET	KM	NCUT	1-SPEC	PNMF	NSC	ONMF	LSD	PLSI	DCD1	DCD1.2	DCD2	DCD5
ORL	0.85	0.90	0.92	0.89	0.90	0.90	0.90	0.90	0.83	0.90	0.91	0.91
MED	0.55	0.57	0.52	0.56	0.55	0.55	0.55	0.55	0.56	0.57	0.57	0.57
VOWEL	0.43	0.40	0.38	0.39	0.36	0.37	0.39	0.39	0.28	0.39	0.38	0.41
COIL20	0.77	0.79	0.77	0.75	0.79	0.79	0.79	0.79	0.74	0.80	0.80	0.80
SEMION	0.57	0.61	0.62	0.58	0.62	0.58	0.58	0.58	0.61	0.61	0.61	0.61
FAULTS	0.10	0.08	0.09	0.09	0.11	0.11	0.11	0.11	0.02	0.08	0.11	0.11
SEGMENT	0.58	0.55	0.58	0.43	0.48	0.43	0.49	0.49	0.07	0.55	0.53	0.58
CORA	0.34	0.16	0.14	0.14	0.13	0.17	0.17	0.17	0.15	0.20	0.24	0.26
CITFSEER	0.34	0.10	0.12	0.07	0.07	0.07	0.07	0.07	0.10	0.18	0.21	0.21
7SECTORS	0.17	0.04	0.05	0.05	0.04	0.04	0.04	0.04	0.04	0.08	0.13	0.11
OPTDIGITS	0.70	0.72	0.80	0.67	0.68	0.67	0.67	0.67	0.69	0.69	0.69	0.69
SVMGUIDE1	0.31	0.35	0.65	0.27	0.27	0.27	0.27	0.27	0.02	0.40	0.39	0.59
ZIP	0.40	0.78	0.79	0.54	0.67	0.65	0.64	0.64	0.21	0.78	0.79	0.81
USPS	0.62	0.77	0.80	0.60	0.75	0.71	0.66	0.66	0.46	0.76	0.77	0.81
PENNDIGITS	0.68	0.81	0.78	0.78	0.78	0.78	0.78	0.78	0.10	0.81	0.83	0.86
PROTEIN	0.00	0.01	0.01	0.02	0.01	0.01	0.01	0.01	0.02	0.01	0.04	0.04
20NEWS	0.05	0.54	0.52	0.36	0.36	0.34	0.34	0.34	0.14	0.45	0.44	0.45
LEFT-REC	0.35	0.38	0.26	0.43	0.43	0.43	0.42	0.42	0.18	0.36	0.42	0.46
MINST	0.51	0.81	0.89	0.59	0.82	0.72	0.59	0.59	0.48	0.80	0.92	0.93

Table 2: The complete clustering results of the first group of experiments.

Table 4: Clustering performance comparison of DCD using *heterogeneous co-initialization* with three ensemble clustering methods. Rows are ordered by dataset sizes. Boldface numbers indicate the best. The ensemble bases are created by the classical k -means algorithm.

DATASET	Purity				NMI			
	BEST	CO	CTS	DCD	BEST	CO	CTS	DCD
ORL	0.77	0.75	0.76	0.83	0.89	0.88	0.88	0.91
MED	0.60	0.64	0.56	0.58	0.58	0.58	0.59	0.58
VOWEL	0.40	0.33	0.40	0.40	0.43	0.34	0.44	0.40
COIL20	0.74	0.63	0.56	0.70	0.80	0.77	0.74	0.80
SEMEION	0.69	0.58	0.62	0.77	0.59	0.54	0.57	0.68
FAULTS	0.42	0.42	0.42	0.44	0.10	0.10	0.10	0.11
SEGMENT	0.59	0.57	0.57	0.65	0.58	0.57	0.56	0.58
CORA	0.60	0.57	0.54	0.55	0.39	0.36	0.37	0.25
CITSEER	0.67	0.68	0.68	0.48	0.38	0.39	0.39	0.21
7SECTORS	0.28	0.30	0.28	0.35	0.07	0.09	0.07	0.11
OPTDIGITS	0.80	0.65	0.74	0.85	0.75	0.70	0.72	0.82
SVMGUIDE1	0.71	0.71	0.71	0.91	0.31	0.31	0.31	0.59
ZIP	0.52	0.41	0.51	0.84	0.42	0.37	0.44	0.81
USPS	0.73	0.73	0.69	0.85	0.61	0.61	0.61	0.81
PENDIGITS	0.77	0.57	0.68	0.89	0.69	0.62	0.68	0.86
PROTEIN	0.46	0.46	0.46	0.50	0.01	0.01	0.01	0.04
20NEWS	0.46	0.42	0.38	0.50	0.58	0.52	0.50	0.45
LET-REC	0.28	0.23	0.27	0.38	0.36	0.29	0.34	0.46
MNIST	0.58	0.48	0.52	0.98	0.49	0.38	0.43	0.93