

# A Decision Criterion for the Optimal Number of Clusters in Hierarchical Clustering

Yunjae Jung ([yunjae@cs.umn.edu](mailto:yunjae@cs.umn.edu))<sup>\*</sup>

*Qwest Communications, 600 Stinson Blvd., Minneapolis, MN 55413*

Haesun Park ([hpark@cs.umn.edu](mailto:hpark@cs.umn.edu))<sup>†</sup>

*Department of Computer Science and Engineering University of Minnesota,  
Minneapolis, MN 55455 and Korea Institute for Advanced Study 207-43  
Cheongryangri-dong, Dongdaemun-gu Seoul 130-012, KOREA*<sup>‡</sup>

Ding-Zhu Du ([dzd@cs.umn.edu](mailto:dzd@cs.umn.edu))

*Department of Computer Science and Engineering University of Minnesota,  
Minneapolis, MN 55455*

Barry L. Drake ([bl Drake@yahoo.com](mailto:bl Drake@yahoo.com))

*CDT, Inc., Minneapolis, MN 55454*

Jan. 21, 2002

**Abstract.** Clustering has been widely used to partition data into groups so that the degree of association is high among members of the same group and low among members of different groups. Though many effective and efficient clustering algorithms have been developed and deployed, most of them still suffer from the lack of automatic or online decision for optimal number of clusters.

In this paper, we define clustering gain as a measure for clustering optimality, which is based on the squared error sum as a clustering algorithm proceeds. When the measure is applied to a hierarchical clustering algorithm, an optimal number of clusters can be found. Our clustering measure shows good performance producing intuitively reasonable clustering configurations in Euclidean space according to the evidence from experimental results. Furthermore, the measure can be utilized to estimate the desired number of clusters for partitional clustering methods as well. Therefore, the clustering gain measure provides a promising technique for achieving a higher level of quality for a wide range of clustering methods.

## 1. Introduction

Clustering refers to the process of grouping patterns so that the patterns are similar within each group and remote between different groups [1]. The distribution of groups can be defined as a cluster configuration.

---

<sup>‡</sup> A part of this work was carried out while H. Park was visiting the Korea Institute for Advanced Study, Seoul, Korea, for her sabbatical leave, from September 2001 to July 2002.

<sup>\*</sup> The work of this author was supported in part by the National Science Foundation grant CCR-9901992

<sup>†</sup> The work of this author was supported in part by the National Science Foundation grant CCR-9901992



The cluster configuration is valid if clusters cannot reasonably occur by chance or as a beneficial artifact of a clustering algorithm [2]. An optimal cluster configuration is defined as an outcome of all possible combinations of groupings, which presents a set of the most “meaningful” associations. Even if the definition of clustering is that simple, evaluation of clustering performance is well known as a fundamental but difficult problem. One reason is that clustering should be performed without *a priori* understanding of the internal structure of the data. In addition, it is impossible to determine which distribution of clusters is best given certain input patterns without an objective measure for clustering optimality. Thus, there have been many attempts to formulate a measure of optimal clustering in the past. However, only a small number of independent clustering criteria can be understood both mathematically and intuitively [12, 2]. Consequently, the hundreds of criterion functions proposed in the literature are related and the same criterion appears in several disguises [3, 4, 5, 6, 7].

Even though an objective measure is given, the difficulty of optimal clustering stems from the astounding number of possible combinations of cluster configurations [8]. The number of ways of generating  $k$  clusters from  $n$  patterns is a Stirling number of the second kind [9, 10, 11]:

$$\mathcal{S}_n^{(k)} = \frac{1}{k!} \sum_{i=1}^k (-1)^{(k-i)} \binom{k}{i} i^n.$$

In particular, the huge volume of data and the potentially high dimensionality of the patterns increase the difficulty of achieving a measure for optimal clustering. Furthermore, it is hard to select a criterion that translates into an intuitive notion of a “cluster” from a reasonable mathematical formula [12]. Feature selection before clustering and cluster labeling after clustering are also challenging problems. As a result, many clustering algorithms to date have been heuristic or ad hoc [12, 2].

Since no ideal solution to the optimal clustering problem has existed from early clustering research [13], recently proposed algorithms have focused mostly on efficiency [14, 15, 16, 17, 18, 19] and scalability [20, 21, 22, 23, 24, 25] to reduce the computational cost and increase processing capability, respectively. It may be possible to produce a cluster configuration very quickly and process huge amounts of data at once. However, often there is no guarantee of achieving an optimal or close-to-optimal clustering configuration.

We propose a method to measure clustering optimality quantitatively with a purpose to use it to determine an optimal number of clusters in various clustering algorithms. The method has been designed

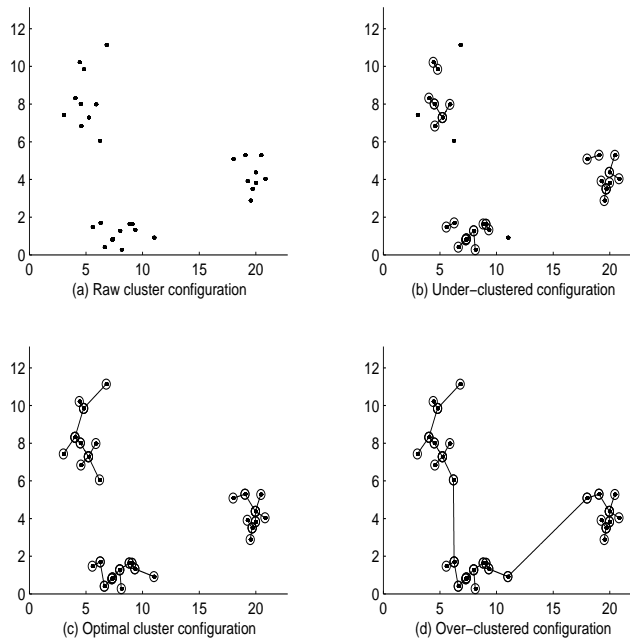


Figure 1. A series of snapshots of clustering configurations in Euclidean distance.

based on the assumption that the optimal cluster configuration can be recognized only by the intuitive and subjective interpretation of a human. Since intuitive validation of clustering optimality can be maximized in two dimensional feature space, it is useful to consider two dimensional Euclidean space for the sake of an objective decision as depicted in Figure 1.

In order to quantify clustering optimality, we introduce clustering gain, which has been designed to have a maximum value when intra-cluster similarity is maximized and inter-cluster similarity is minimized [26, 8]. Thus, the optimal cluster configuration can be identified by the maximum of the clustering gain curve. This measure can be directly used to explore an optimal configuration for all hierarchical clustering algorithms as they proceed. The measure can also be useful for performance comparison among clustering algorithms since the clustering performance is also measured by clustering gain.

Since discovering all possible combinations of cluster configuration is computationally prohibitive [8, 12], most partitional clustering algorithms are dependent on users to determine the target number of clusters. We show how the desired number of clusters can be estimated based on the data using the clustering gain measure. The best cluster configuration will be one which can be produced by any specific

hierarchical clustering algorithm. Consequently, the measure is used to estimate the desired number of clusters in partitional clustering algorithms.

According to the experimental results, most commonly used hierarchical clustering algorithms are able to produce intuitively reasonably clustered configurations using our clustering measure when the input patterns are distributed in a well isolated fashion. Moreover, the desired number of clusters for partitional clustering methods, e.g.  $k$ -means, has been successfully estimated experimentally.

The rest of this paper is organized as follows. In Section 2, some background information on previous work is presented to derive optimal clustering measures. The design scheme of the optimal clustering gain measure is discussed in Section 3. Section 4 discusses how the proposed method can be used to evaluate the performance of clustering algorithms. Finally, we discuss how to estimate the optimal number of clusters for partitional clustering algorithms using our new scheme in Section 5.

## 2. Optimal Clustering

Stopping criteria for optimal clustering have been the topic of considerable past research effort [27]. Specifically, deciding the optimal level of a dendrogram and estimating the number of target clusters remains as a challenging and fundamental problem. For hierarchical agglomerative clustering, some decision rules have been provided by Milligan and Cooper [28] to determine the appropriate level of the dendrogram [29, 30] for optimal clustering. In addition, Milligan compared and described objective functions for optimal agglomerative clustering. However, the functions are based on strong assumptions, heuristics, and experimental dependency. Recently, a stopping rule for the hierarchical divisive clustering method has been suggested in the Principal Direction Divisive Partitioning (PDDP) algorithm [31]. While the PDDP algorithm is proceeding, a dynamic threshold based on a so-called centroid scatter value is calculated. The rule is to stop partitioning when the centroid scatter value exceeds the maximum cluster scatter value at any particular point. This approach relies on experimental observations. In general, currently used stopping criteria for hierarchical clustering methods are based on predefined thresholds including the number of iterations, the number of clusters, average dissimilarity [32] within a cluster, maximum distance between patterns, and relative inter-connectivity and relative closeness [33].

For non-hierarchical partitioning algorithms, Dubes [27] provided a separation index:

$$S(k) = \frac{|f(k+1, k) - f(k, k-1)|}{1 + |f(k+1, k)f(k, k-1)|}, \quad (1)$$

where

$$f(k+1, k) = MH(k+1) - MNH(k).$$

The value  $MH$  is the point serial correlation coefficient between the matrix of Euclidean distances for patterns and a “model” matrix, and  $k$  is the number of clusters. The model matrix sets the distance between two patterns to be the distances between the centers of clusters to which the patterns belong. A stopping rule is adopted to search for a significant knee in the curve of  $MH(k)$  as  $k$  varies from  $k_{max}$  to 2 where  $k_{max}$  is the maximum possible number of clusters. However, a threshold that distinguishes the knee from other anomalies is difficult to determine. In addition, the rule is not able to avoid premature stopping, i.e., the convergence to local minima problem. Similarly, Boulder and Odell [34] introduced a cluster separation measure

$$R_{ij} \equiv \frac{S_i + S_j}{M_{ij}}, \quad (2)$$

where  $S_i$  is a dispersion measure of cluster  $i$  such as the squared error sum and  $M_{ij}$  is the distance between two centroids. The separation measure will be that which minimizes the average similarity as follows

$$\bar{R} \equiv \frac{1}{n} \sum_{i=1}^k n_i R_i, \quad (3)$$

where  $R_i \equiv$  maximum of  $R_{ij}, i \neq j$ ,  $n_i$  is the number of patterns in cluster  $i$ , and  $n$  is the total number of the patterns to be clustered. However, a unique optimal clustering condition can not be detected by the separation measure. According to the experimental results, there are many minimum points in the hierarchical system [34]. Furthermore, there is no theoretical basis for the feasibility of the measure and no reasonable separation measure for partitioning clustering seems to exist at present [19].

### 3. Design of A Measure for Optimal Clustering

#### 3.1. CLUSTERING BALANCE

The clustering problem is to partition the given input patterns into a specific number of groups (clusters) so that the intra-cluster similarity is maximized and the inter-cluster similarity is minimized in a particular metric space [26, 8]. Throughout the paper, we will use the following notations. Pattern  $i$  is a feature vector in an  $m$  dimensional space, denoted as

$$p_i = [p_{i1}, p_{i2}, \dots, p_{im}]^T,$$

and a cluster  $C_j$  is a set of patterns grouped together by a clustering algorithm and expressed by

$$C_j = \{p_1^{(j)}, p_2^{(j)}, \dots, p_{n_j}^{(j)}\},$$

where  $n_j$  is the number of patterns in cluster  $C_j$ . We will assume that there are total of  $n$  vectors to be clustered and the total number of the clusters is  $k$ . Accordingly,  $\sum_{i=1}^k n_i = n$ . In addition,  $p_0^{(j)}$  denotes the centroid of the cluster  $j$ , which is defined as

$$p_0^{(j)} = \frac{1}{n_j} \sum_{i=1}^{n_j} p_i^{(j)}.$$

The centroid is a mean vector of the cluster and provides a compressed representation of the cluster in a *simpler* form. The centroid is often used for cluster data compression.

Cluster configuration is a random variable whose possible outcome is a particular assignment of input pattern sets. The problem of optimal clustering is to find a cluster configuration that is optimized according to some evaluation criterion. However, as mentioned before, the number of ways of clustering  $n$  observations into  $k$  groups is enormously large [35]. In fact, a combinatorial search of the set of possible configurations for optimal clustering is clearly computationally prohibitive [27] and, is NP-complete [36, 37]. Accordingly, currently used agglomerative clustering algorithms take an approximation approach by merging more similar patterns prior to grouping less similar patterns to construct a cluster hierarchy. A measure of the similarity between two patterns drawn from the same feature space plays an essential role in these clustering procedures [12, 35].

The most popular metric for measuring similarity between patterns is the Euclidean distance since it is more intuitive and applicable, especially in two dimensional feature space [13, 2]. The most intuitive

and frequently used criterion function in clustering techniques is the squared error criterion which is the sum of squared distances from the centroid of a group to every pattern in the group [38, 13, 39, 14, 40, 41, 42, 43, 19] which can be expressed using the Euclidean distance [2, 44].

The intra-cluster error sum  $\Lambda$  is defined by the squared error  $e$  as

$$\Lambda = \sum_{j=1}^k \sum_{i=1}^{n_j} e(p_i^{(j)}, p_0^{(j)}).$$

which, using the Euclidean distance, can be denoted as

$$\Lambda = \sum_{j=1}^k \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0^{(j)}\|_2^2. \quad (4)$$

It is also called the within-group error sum [12]. Ward used the error sum of squared to quantify the loss of information by grouping [38].

The inter-cluster error sum takes into account error sums between clusters by considering the collection of cluster centroids to be a global pattern set, which also has a global centroid. The inter-cluster error sum, in case of Euclidean space is defined as

$$\Gamma = \sum_{j=1}^k e(p_0^{(j)}, p_0) = \sum_{j=1}^k \|p_0^{(j)} - p_0\|_2^2, \quad (5)$$

where  $p_0$  is the global centroid defined as

$$p_0 = \frac{1}{n} \sum_{i=1}^n p_i.$$

Now, we present some characteristics of these two conflicting error sums, to be utilized in designing a measure for optimal cluster configuration as well as a stopping criterion in hierarchical clustering algorithm. We will assume that the hierarchical algorithm we are considering is agglomerative. In case of divisive algorithms, an analogous but opposite trends can easily be proved. We can assume that in the initial state of any agglomerative clustering algorithm, each pattern is the only pattern in its own cluster. It is clear that the singleton clusters have no contribution to the intra-cluster error sum  $\Lambda$ , and the minimum value that  $\Lambda$  can take is zero. On the other hand,  $\Lambda$  is maximized when there is only one cluster that contains all patterns. More interesting fact is that while clustering process proceeds, the value of  $\Lambda$  cannot decrease. Suppose two clusters  $C_i$  and  $C_j$  are merged in a step of agglomerative clustering. Let the new cluster  $C_{ij}$  be the cluster obtained by merging

$C_i$  and  $C_j$ . Then, the centroid  $c_{ij}$  of the new cluster  $C_{ij}$  is

$$c_{ij} = \frac{n_i p_0^{(i)} + n_j p_0^{(j)}}{n_i + n_j}.$$

Let  $\Lambda_b$  and  $\Lambda_a$  be intra-cluster error sums of the items that belong to the clusters  $C_i$  and  $C_j$  only, before and after merging, respectively.

Then

$$\Lambda_b = \sum_{l=1}^{n_i} \|p_l^{(i)} - p_0^{(i)}\|_2^2 + \sum_{l=1}^{n_j} \|p_l^{(j)} - p_0^{(j)}\|_2^2,$$

and

$$\Lambda_a = \sum_{l=1}^{n_i} \|p_l^{(i)} - c_{ij}\|_2^2 + \sum_{l=1}^{n_j} \|p_l^{(j)} - c_{ij}\|_2^2.$$

Since, there is no split of a cluster in a path of agglomerative clustering, intra-cluster error sum would be nondecreasing as the clustering proceeds if  $\Lambda_a - \Lambda_b \geq 0$ . We have

$$\begin{aligned} \Lambda_a - \Lambda_b &= \sum_{l=1}^{n_i} \|p_l^{(i)}\|_2^2 - 2c_{ij}^T \sum_{l=1}^{n_i} p_l^{(i)} + n_i c_{ij}^T c_{ij} + \sum_{l=1}^{n_j} \|p_l^{(j)}\|_2^2 \\ &\quad - 2c_{ij}^T \sum_{l=1}^{n_j} p_l^{(j)} + n_j c_{ij}^T c_{ij} \\ &\quad - \left[ \sum_{l=1}^{n_i} \|p_l^{(i)}\|_2^2 - 2(p_0^{(i)})^T \sum_{l=1}^{n_i} p_l^{(i)} + n_i \|p_0^{(i)}\|_2^2 + \sum_{l=1}^{n_j} \|p_l^{(j)}\|_2^2 \right. \\ &\quad \left. - 2(p_0^{(j)})^T \sum_{l=1}^{n_j} p_l^{(j)} + n_j \|p_0^{(j)}\|_2^2 \right]. \end{aligned}$$

Using  $\sum_{l=1}^{n_i} p_l^{(i)} = n_i p_0^{(i)}$  and  $\sum_{l=1}^{n_j} p_l^{(j)} = n_j p_0^{(j)}$ , we have the desired result

$$\begin{aligned} \Lambda_a - \Lambda_b &= 2n_i \|p_0^{(i)}\|_2^2 - n_i \|p_0^{(i)}\|_2^2 - 2n_i (p_0^{(i)})^T c_{ij} + n_i \|c_{ij}\|_2^2 \\ &\quad + 2n_j \|p_0^{(j)}\|_2^2 - n_j \|p_0^{(j)}\|_2^2 - 2n_j (p_0^{(j)})^T c_{ij} + n_j \|c_{ij}\|_2^2 \\ &= n_i \|p_0^{(i)} - c_{ij}\|_2^2 + n_j \|p_0^{(j)} - c_{ij}\|_2^2 \geq 0. \end{aligned}$$

Similarly, the inter-cluster error sum satisfies the following characteristics which show the opposite trend to that of the intra-cluster error sum. Note that the global centroid  $p_0$  does not change throughout the clustering process. The inter-cluster error-sum  $\Gamma$  is maximized when there are  $n$  singleton clusters, which occurs at the beginning of



clustering. Then  $\Gamma$  is minimized when all  $n$  patterns belong to one cluster at the end of clustering. It is easy to show that the value of  $\Gamma$  is nonincreasing as the clustering proceeds using the triangular property of the  $L_2$  norm, the Euclidean distance.

Our design scheme is based on the fact that intra-cluster similarity is nondecreasing and inter-cluster error sum is nonincreasing as the agglomerative clustering algorithm proceeds. When the clustering algorithm is divisive, the trend is the other way around, which is that intra-cluster error sum is nonincreasing and inter-cluster similarity is nondecreasing as the divisive clustering algorithm proceeds.

We transformed the optimal clustering problem into a problem for finding the point where the two similarities are balanced by representing these similarities by the squared error sum in Euclidean space. We define the *clustering balance* as

$$\mathcal{E}(\chi) = \alpha\Lambda + (1 - \alpha)\Gamma, \quad (6)$$

where  $\Lambda$  and  $\Gamma$  denote intra-cluster and inter-cluster error sums for a specific clustering configuration  $\chi$ , respectively, and  $0 \leq \alpha \leq 1$  is a scalar that determines the weight between these two sums. The clustering balance  $\mathcal{E}(\chi)$  has been formulated with the idea that intuitively optimal clustering is achieved when the error sums have reached equilibrium. We will concentrate on the special case for  $\alpha = 1/2$  which provides an even balance and accordingly assume that

$$\mathcal{E}(\chi) = \Lambda + \Gamma. \quad (7)$$

Thus clustering behavior can be interpreted as a procedure seeking the global minimum of clustering balance. With the definitions of clustering balance based on the error sums, what follows will use the trade-off between inter-cluster and intra-cluster balance to define a measure for the optimal clustering configuration.

### 3.2. CLUSTERING GAIN

The clustering balance can be computed in each step of a hierarchical clustering algorithm to determine the optimal number of clusters. However, a major disadvantage is the high computational cost of computing clustering balance. In this section, we introduce *clustering gain* which has an interesting relation to cluster balance. In addition, clustering gain is cheap to compute. Therefore, it can be computed in each step of clustering process to determine the optimal number of clusters without increasing the computational complexity.

Clustering gain  $\Delta_j$  for  $C_j$  is defined as the difference between the decreased inter-cluster error sum  $\gamma_j$  compared to the initial stage and

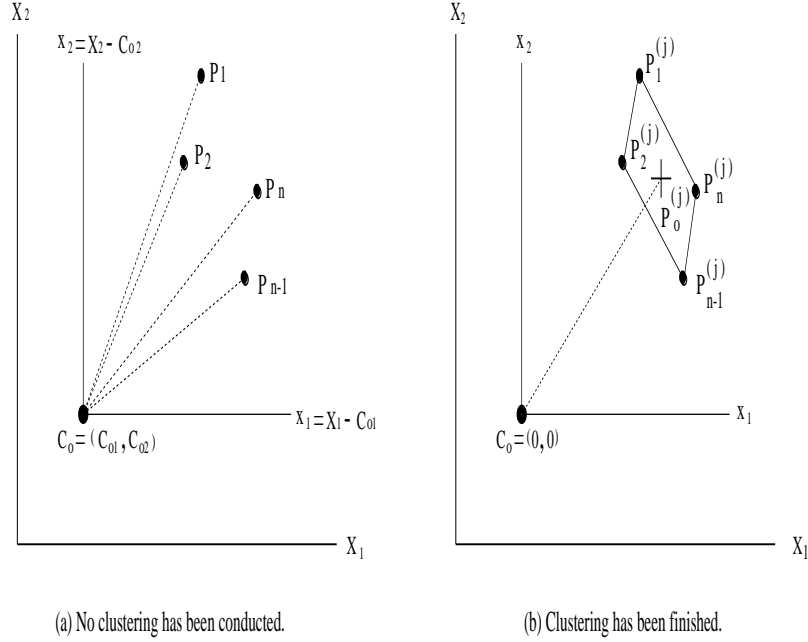


Figure 2. Clustering gain defined by difference between error sums. (a) Initial configuration. (b) Final configuration of cluster  $C_j$ .

the increased intra-cluster error sum  $\lambda_j$  compared to the initial stage. Specifically, the gain is defined as

$$\Delta_j = \gamma_j - \lambda_j.$$

In the above equation, an equal weighting factor, one, has been assigned to both error sums. Clustering gain is graphically illustrated in Figure 2 using cluster configurations. In particular, the number of patterns of the final configuration of cluster  $C_j$  can vary from 1 to  $n$ .

The decreased portion of the inter-cluster error sum compared to the initial stage is denoted by

$$\begin{aligned} \gamma_j &= \sum_{i=1}^{n_j} e(p_i^{(j)}, p_0) - e(p_0^{(j)}, p_0) \\ &= \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0\|_2^2 - \|p_0^{(j)} - p_0\|_2^2 \end{aligned}$$

In addition, the increased portion of the intra-cluster error sum compared to the initial stage is defined by

$$\lambda_j = \sum_{i=1}^{n_j} e(p_i^{(j)}, p_0^{(j)}) = \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0^{(j)}\|_2^2. \quad (8)$$

Expanding the gain for cluster  $C_j$  gives

$$\begin{aligned}\Delta_j &= \gamma_j - \lambda_j \\ &= \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0\|_2^2 - \|p_0^{(j)} - p_0\|_2^2 - \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0^{(j)}\|_2^2 \\ &= (n_j - 1) \|p_0 - p_0^{(j)}\|_2^2\end{aligned}$$

since  $\sum_{i=1}^{n_j} p_i^{(j)} = p_0^{(j)} n_j$ . Therefore, the total clustering gain can be computed from

$$\Delta = \sum_{j=1}^k (n_j - 1) \|p_0 - p_0^{(j)}\|_2^2. \quad (9)$$

We would like to emphasize that this clustering gain is very cheap to compute since it involves only the centroids and the global centroid, and not the individual data items. The clustering gain  $\Delta_j$  is always greater than or equal to zero. Eventually, the clustering gain will be positive, assuming the initial clustering configuration is not optimal.

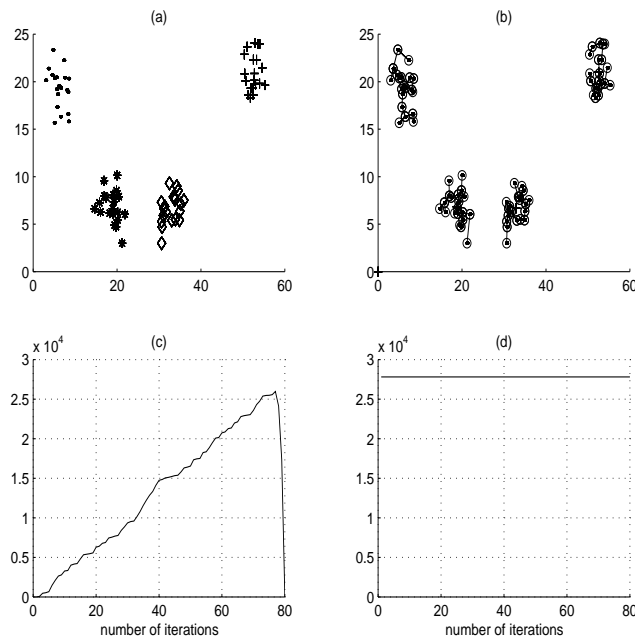
Apparent from Figure 3 is the fact that the optimal clustering configuration discovered by a hierarchical clustering algorithm has maximum clustering gain. Since clustering gain is minimum at the initial and final clustering stages, an optimal configuration should be found during the middle of the clustering procedure. In order to determine the maximum clustering gain during the middle of the clustering procedure, we propose the *clustering gain* as an effectiveness criterion. Note that clustering gain is analogous to the  $E$  value suggested by Jardine and Rijsbergen [45], for clustering effectiveness.

It is interesting to note that the sum of clustering balance and clustering gain is a constant for a given data set since

$$\begin{aligned}\Omega &= \mathcal{E} + \Delta = \Lambda + \Gamma + \Delta \\ &= \sum_{j=1}^k \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0^{(j)}\|_2^2 + \sum_{j=1}^k \|p_0^{(j)} - p_0\|_2^2 \\ &+ \sum_{j=1}^k \left( \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0\|_2^2 - \|p_0^{(j)} - p_0\|_2^2 - \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0^{(j)}\|_2^2 \right) \\ &= \sum_{j=1}^k \sum_{i=1}^{n_j} \|p_i^{(j)} - p_0\|_2^2\end{aligned}$$

which is determined completely based on the data, and not changed by the clustering result. Thus clustering balance can be alternatively expressed using clustering gain as

$$\mathcal{E} = \Lambda + \Gamma = \Omega - \Delta,$$



*Figure 3.* Clustering gain: the opposite concept of clustering balance. (a) Initial configuration of patterns. (b) Intuitively well clustered configuration captured when clustering gain is maximized. (c) Clustering gain (d) The sum of clustering balance and clustering gain.

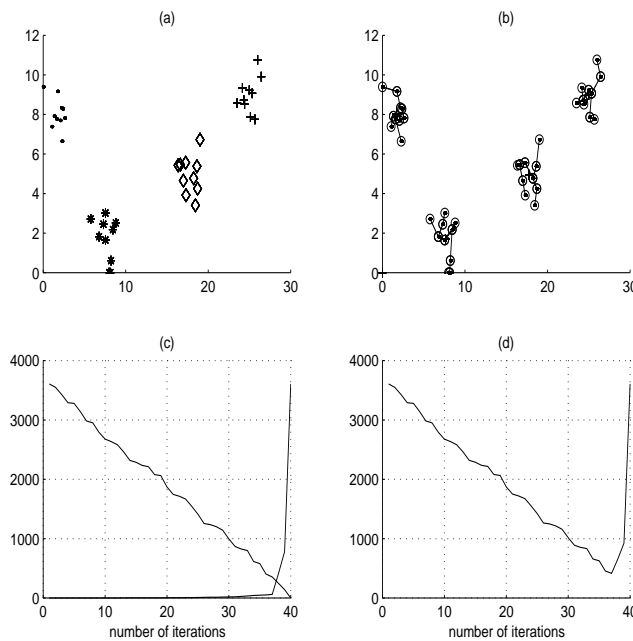
where  $0 \leq \Lambda, \Gamma, \Delta \leq \Omega$ . Now, we are able to find an optimal cluster configuration by tracing clustering gain instead of clustering balance. For visual demonstration, clustering gain  $\Delta$  and the constant  $\Omega$  are compared in Figure 3.

In Tracking Algorithm, we summarize how we can obtain the optimal cluster configuration in a given hierarchical agglomerative clustering algorithm while keeping track of the cluster gain value. Note that we need to keep track of the clustering gain  $\Delta(\chi)$  since the global maximum value of clustering gain can be discovered only after the clustering is completed.

Application of this method to hierarchical divisive clustering method is straightforward. To demonstrate the performance of Tracking Algorithm, the optimal configuration detected by the complete-link is visualized in Figure 4. Clearly, the configuration containing the lowest value of clustering balance coincides with the optimal configuration produced by human being. Given the same input patterns, the same optimal configuration has been obtained by popular agglomerative clustering algorithms including the single-link, the average-link and Ward's method.

**Tracking Algorithm:**

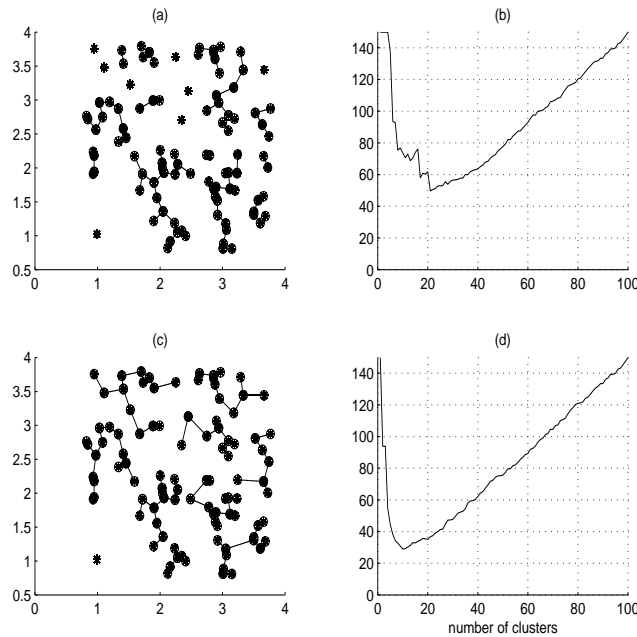
- 1 Choose a hierarchical agglomerative clustering method(HACM)
- 2 Do while clustering is not complete in HACM
- 4 merge two clusters according to the fusion rule of the HACM
- 5 keep track of maximum value for  $\Delta(\chi)$  and save  $\chi$
- 6 end while
- 7 recover the optimal configuration  $\chi$



*Figure 4.* Visual demonstration of optimal configuration discovered by Tracking Algorithm using the complete-link. (a) Initial configuration. (b) Optimal configuration. (c) Intra-cluster versus inter-cluster error sum. (d) Clustering balance.

#### 4. Performance Evaluation of Clustering Algorithms

Given a hierarchical clustering algorithm, either clustering balance or clustering gain can be used to find optimal configuration. Since these measures represent clustering optimality in an absolute value, they also can be used to compare various clustering algorithms in terms of clustering performance. To give an example we estimated the practical distribution in high dimensionality using the method proposed by Bennett [46]. The method is based on the observation that the variance



*Figure 5.* Comparison of the single-link and the average-link. (a) Optimal configuration found by single-link. (b) Clustering balance versus the number of clusters in single-link. (c) Optimal configuration found by average-link. (d) Clustering balance versus the number of clusters in average-link.

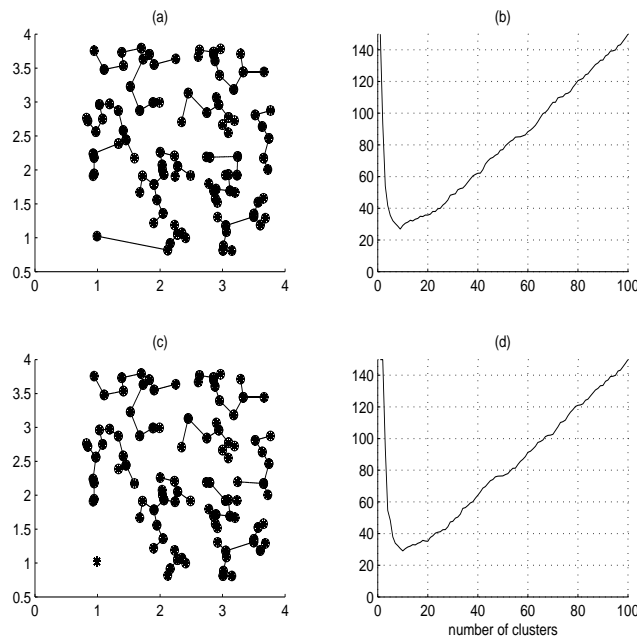
of the distance between patterns chosen randomly in a hyper-sphere is inversely proportional to the dimensionality  $m$  as follows.

$$m \times \text{variance}(\text{a pair of patterns}) \simeq \text{constance}.$$

According to the equation, the higher the dimensionality is the smaller the variance is. Thus clustering in high dimensional space can be simulated in two dimensional space if patterns are randomly and uniformly distributed. Typically used agglomerative clustering algorithms have been applied to Tracking Algorithm, and then their optimal configurations and clustering balances are represented in Figure 5 and Figure 6.

According to the results, the complete-link produces the configuration with the lowest clustering balance. In the experiment, the complete-link outperformed other three clustering algorithms since it produces the best configuration given the same input patterns.

To extend our approach to practical domain, we conducted a simple experiment with practical document vectors. The documents have been downloaded from MEDLINE on-line library. They are divided into eight cancer categories including breast, colon, weightless, glycolic, heart attack, oral, prostate and tooth-decay. Each category con-



*Figure 6.* Comparison of the complete-link and Ward's method. (a) Optimal configuration found by the complete-link. (b) Clustering balance versus the number of clusters in the complete-link. (c) Optimal configuration found by Ward's method. (d) Clustering balance versus the number of clusters in Ward's method.

tains 500 documents. After filtering the document set using the stoplist of SMART IR system and the stemming algorithm proposed by Porter [48], we applied Tracking Algorithm to the combination of Colon and Tooth categories. The results in Euclidean space are graphically illustrated in Figure 7. According to the results, optimal cluster configurations can be found by our measure in Euclidean space.

## 5. Estimation of the Optimal Number of Clusters for Partitional Clustering Algorithms

A major problem accompanying the use of a partitional clustering algorithm is the choice of the number of desired output clusters [27] and order dependency [13]. The sensitivity to the selection of the initial partition is a problem in most partitional algorithms [2] since the partitional clustering algorithms may converge to a local minimum of the criterion function if the initial partition is not properly chosen [2]. For handling large data set and avoiding such sensitivity, many efficient algorithms have been proposed including CLARANS [19], BIRCH [47], CURE [23], and ROCK [24]. Many clustering algorithms have efficiency

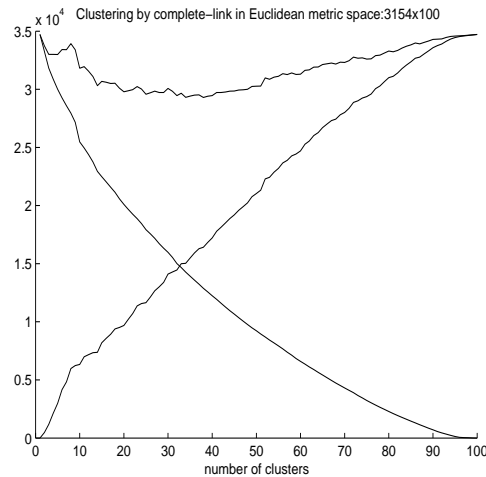


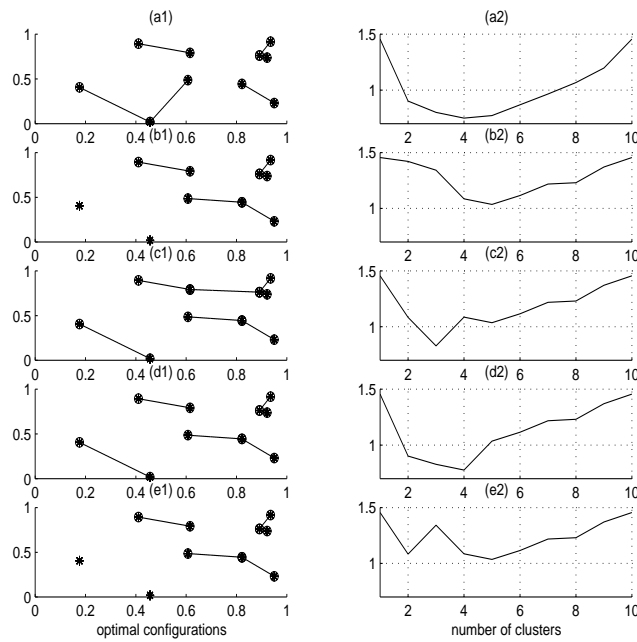
Figure 7. Clustering balance, intra-cluster and inter-cluster squared error sums in Euclidean space.

and capacity, but most partitional clustering algorithm depend on users to determine the desired number of clusters.

Even though exhaustive enumeration of all possible assignments is not computationally feasible even for small numbers of patterns [12], we can generate all possible cases with very small number of patterns such as ten. In Figure 8, optimal configurations and their clustering balances are compared with respect to ideal configuration. For each possible number of clusters, the lowest balance is presented in the part (a) of the figure. According to the experimental results, the complete-link discovered the optimal configuration as closely as the ideal. However, it is risky to generalize this result so that the complete-link is superior to all other algorithms for all input patterns since clustering performance may change in accordance with the distribution of the input patterns. A hierarchical clustering algorithm is considered as the best if it produces the lowest clustering balance given particular data patterns.

As we previously described, the best configuration can be selected among optimal configurations produced by hierarchical clustering algorithms. Consequently, the desired number of clusters can be estimated from the best configuration. Also the centroids of the best configuration can be fed to partitional clustering algorithms to avoid random initial assignments for centroids. The basic assumption of this approach is that the best cluster configuration, the winner among configurations produced by hierarchical clustering algorithms, will be an approximation of the ideal cluster configuration for partitional clustering.





*Figure 8.* Comparison of currently used agglomerative clustering algorithms to the optimal clustering by using cluster configurations and clustering balances. (a) The optimal (b) The single-link (c) The average-link (d) The complete-link (e) Ward's method.

It is clear that the estimated number is not the true value. However, the estimated number can contribute to decision of the range of the true number of optimal clusters. To verify this assumption experimentally, we applied  $k$ -means algorithm with all possible number of clusters. The averaged clustering balance produced by the  $k$ -means is depicted in Figure 9 along with the number of clusters. In this experiment, the clustering balance is the average of five trials. According to the experimental results, the desired number of clusters for the given distribution is nine. When we apply Tracking Algorithm to four popular algorithms including single-link, average-link, complete-link and Ward's method, corresponding optimal configurations are found as in Table I. Surprisingly, the number of clusters produced by the complete-link is equivalent to the desired number of clusters obtained by  $k$ -means using all possible  $k$  values. This result convincingly illustrates that our clustering measure can be used for partitional algorithms to estimate the number of desired clusters. In addition, more stable configuration and improved performance are demonstrated in Figure 10. When the number of desired clusters and initial centroids are estimated,  $k$ -means is able to converge quickly.

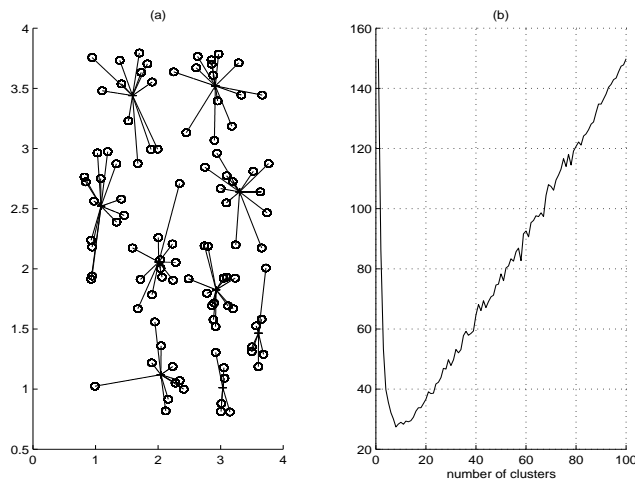


Figure 9. The optimal configuration and clustering balance traced by  $k$ -means. (a) Optimal configuration. (b) Averaged clustering balance of all possible number of clusters using  $k$ -means.

Table I. Comparison of lowest balances

algorithms	the lowest balance	the highest gain	number of clusters
Single-link	49.8603	100.139651	21
Average-link	29.2123	120.787653	10
Complete-link	27.0904	122.909583	9
Ward's Method	29.2292	120.770844	10

For  $k$ -medoid algorithm, the results are almost the same as  $k$ -means algorithm except some fluctuations of clustering balance before convergence. This result is normal since centroids and medoids are located differently in the same feature space. As a result, the best cluster configuration found by hierarchical clustering algorithms contributes to determining the desired number of clusters and the initial centroids for partitional clustering algorithms.

## 6. Conclusion

Clustering is not a new technique in computer related disciplines. However, a huge demand for clustering technique represented by a variety

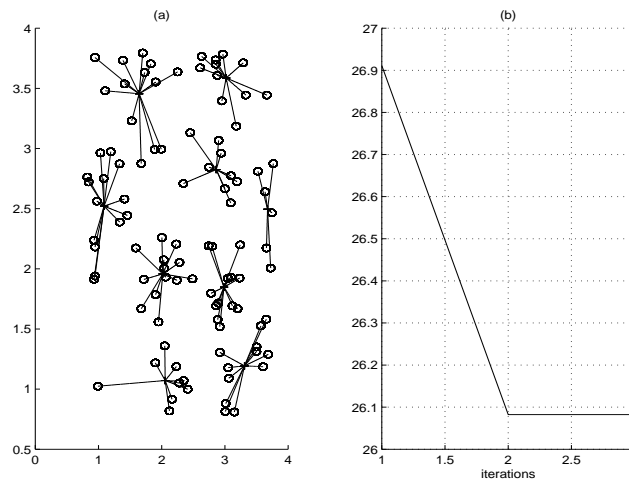


Figure 10. Additional clustering optimization and quick convergence. (a) Optimal configuration found by a hierarchical clustering algorithm. (b) Improved clustering performance.

of clustering applications demonstrates its importance. Recently, much effort has been presented to achieve clustering efficiency and scalability. In this paper, we proposed a measure for optimal clustering. We defined clustering balance using the squared error sums. By searching the compromising point between intra-cluster and inter-cluster error sums, we are able to detect the optimal clustering configuration for any hierarchical clustering algorithms. For the purpose of finding an optimal configuration, an agglomerative clustering recovers the cluster configuration with the minimum clustering balance.

Our approach is quite different from other traditional approaches. We evaluate clustering optimality using only internal properties of clusters and successfully achieve intuitive agreement for clustering optimality. In particular, the number of desired clusters and initial centroids can be estimated from the optimal cluster configuration, and provided to non-hierarchical partitional clustering methods. As a result, partitional clustering algorithms are able to converge more quickly and give lower clustering balance than those without our clustering measure.

When it comes to classification, multiple centroids in a class can be found using our clustering measure since each class in turn is a cluster. It is natural to assume that those centroids provide us with more accurate information describing the internal structure of a class than that represented by only one centroid. Therefore, we believe that classification performance will be enhanced if we exploit our clustering measure to find optimal sub-centroids in each class. The basic

rationale of improved classification is that classification performance is contributed by comparing test data to multiple centroids instead of the single centroid. Therefore, we expect that many unknown contributions of our approach will be discovered in various applications of clustering while our clustering measure consistently gives feasible solutions to optimal clustering.

## References

1. E. Gose, R. Johnsonbaugh and S. Jost, *Pattern Recognition & Image Analysis*, Prentice Hall, Upper Saddle River, NJ, 1996
2. A.K. Jain and M.N. Murty and P.J. Flynn, 'Data Clustering: A Review', *ACM Computing Surveys*, **Vol.31 no.3**, pp. 264–323, 1999
3. E. Shaffer, R. Dubes and A.K. Jain, 'Single-link Characteristics of a Mode-seeking Algorithm', *Pattern Recognition*, **Vol.11**, pp. 65–73, 1979
4. J. Kittler, 'A Locally Sensitive Method for Cluster Analysis', *Pattern Recognition*, **Vol.8**, pp. 22–33, 1976
5. C.T. Zahn, 'Graph-theoretical methods for Detecting and Describing Gestalt Clusters', *IEEE Transactions on Computers*, **Vol.20**, pp. 68–86, 1971
6. R. Urquhart, 'Graph Theoretical Clustering based on limited neighborhood sets', *Pattern Recognition*, **Vol.15**, pp. 173–187, 1982
7. K.C. Gowdar and G. Krishna, 'Agglomerative Clustering using the Concept for Multispectral Data', *Pattern Recognition*, **Vol.10**, pp. 105–112, 1978
8. M.R. Anderberg, *Cluster Analysis for Applications*, Academic Press, New York, 1973
9. M. Abramowitz and I.A. Stegun, *Handbook of Mathematical Functions with Formulas, Graphics and Mathematical Tables*, US Govt. Printing Office, Washington, D.C., 1968
10. J.J. Fortier and H. Solomon, *Clustering Procedures*, In *Multivariate Analysis*, P. R. Krishnaiah, editor. Academic Press, Inc., New York, NY, pp. 493–506, 1966
11. R.E. Jensen, 'A Dynamic Programming Algorithm for Cluster Analysis', *Operations Research*, **Vol.17**, pp. 1034–1057, 1967
12. A.K. Jain and R.C. Dubes, *Algorithms for Clustering Data*, Prentice Hall, Englewood Cliffs, NJ, 1988
13. P. Willet, 'Recent Trends in Hierarchic Document Clustering: A Critical Review', *Information Processing and Management*, **Vol.24 no.5**, pp. 577–597, 1988
14. R.C. Jancey, 'Multidimensional Group Analysis', *Austral. J. Botany*, **Vol.14 no.1**, pp. 127–130, 1966
15. J.B. MacQueen, 'Some Methods for Classification and Analysis of Multivariate Observations', In *the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, **Vol.1**, AD 669871, Univ. of California Press, Berkeley, pp. 281–297, 1967
16. R. Sibson:73, 'SLINK: an Optimally Efficient Algorithm for the Single-link Cluster Method', *Computer Journal*, **Vol.16**, pp. 30–34, 1973
17. D. Defays, 'An Efficient Algorithm for a Complete Link Method', *Computer Journal*, **Vol.20**, pp. 364–366, 1977

18. W.H.E. Day and H. Edelsbrunner, 'Efficient Algorithms for Agglomerative Hierarchical Clustering Methods', *Journal of Classification*, **Vol.1 no.1**, pp. 7–24, 1984
19. R.T. Ng and J. Han G., 'Efficient and Effective Clustering Methods for Spatial Data Mining', In *Proceedings of the 20th VLDB Conference*, Santiago, Chile, 1994
20. E.M. Voorhees, 'Implementing Agglomerative hierarchical Clustering Algorithms for use in Document Retrieval', *Information Processing and Management*, **Vol.22**, Pergamon Press Ltd., Great Britain, pp. 465–476, 1986
21. X. Li, 'Parallel Algorithms for Hierarchical Clustering and Cluster Validity', *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **Vol.12 no.11**, pp. 1088–1092, 1990
22. P. Bradley and U. Fayyad and C. Reina, 'Scaling Clusterin Algorithms to Large Databases', *Knowledge Discovery and Data Mining*, 1998
23. S. Guha and R. Rastogi and K. Shim, 'CURE: An Efficient Clustering Algorithm for large Databases', In *ACM-SIGMOD Int. Conf. on Management of Data*, Seattle, WA. USA, pp. 73–84, 1998
24. S. Guha and R. Rastogi and K. Shim, 'ROCK: A Robust Clustering Algorithm for Categorical Attributes', In *the 15th Int. Conf. on Data Eng.*, 1999
25. V. Ganti and R. Ramakrishnan and J. Gehrke, 'Clustering Large Datasets in Arbitrary Metric Spaces', *ACM*, 2000
26. M. Charikar and C. Chekuri and T. Feder and R. Motwani, 'Incremental Clustering and Dynamic Information Retrieval', In *STOC'97*, El Paso, Texas, USA, pp. 153-180, 1997
27. R.C. Dubes, 'How many Clusters are best? – an experiment', *Pattern Recognition*, **Vol.20 no.6**, pp. 645–663, 1987
28. G.W. Milligan and M.C.Cooper, 'An Examination of Procedures for Detecting the number of Clusters in a data set', *Psychometrika*, **Vol.50**, pp. 159–179, 1985
29. J.T. Tou and R.C. Gonzalez, *Pattern Recognition Principles*, Addison-Wesley, Miami, FL, 1974
30. B. Everett, *Cluster Analysis*, Addison-Wesley, New York, 1975
31. D.L. Boley, 'Principal Direction Divisive Partitioning', *Data Mining and Knowledge Discovery*, **Vol.2 no.4**, pp. 325–344, 1998
32. B. Mirkin and I. Muchnik, *Combinatorial Optimization in Clustering*, In *Handbook of Combinatorial Optimization*, D.Z. Du and P.M. Pardalos, Kluwer Academic Publishers, pp. 261–329, 1998
33. G. Karypis and E.S. Han and V. Kumar, 'CHAMELEON: A Hierarchical Clustering Algorithm Using Dynamic Modeling', *IEEE Computer: Special Issue on Data Analysis and Mining*, **Vol. 32 no. 8**, pp 68-75, 1999
34. B.S. Duran and P.L. Odell, *Cluster Analysis: A Survey*, Springer-Verlang, Berlin, 1977
35. E. Diday and J.C. Simon, *Clustering Analysis*, In *Digital Pattern Recognition*, K.S. Fu, Ed., Springer-Verlag, Secaucus, NJ, pp. 47–94, 1976
36. M.R. Garey and D.S. Johnson, *Computers and Intractability: a guide to the theory of NP-completeness*, W.H.Freeman and Company, San Francisco, CA, 1979
37. P. Crescenzi and V. Kann, *A compendium of NP optimization problems*, URL site:<http://www.nada.kth.se/~viggo/problemlist/compendium2>, 1995
38. J.H.Jr. Ward, 'Hierarchical Grouping to Optimize an Objective Function', *Journal of the American Statistical Association*, **Vol.58**, pp. 236–244, 1963
39. E.W. Forgy, 'Cluster Analysis of Multivariate Data: Efficiency Versus Interpretability of Classification', In *Biometric society meetings, Riverside, CA*, Abstract in *Biometrics* **Vol.21 No.3**, pp. 768, 1965

40. G.S. Sebestyen, 'Pattern Recognition by an Adaptive Process of Sample Set Construction', *IRE Trans. on Info. Theory*, **Vol.IT-8**, Sep., 1962
41. J.B. MacQueen, 'Some Methods for Classification and Analysis of Multivariate Observations', In *Wester Management Science Inst.*, University of California, pp. 96, 1966
42. G.H. Ball and D.J. Hall, 'Some Fundamental Concepts and Synthesis Procedures for Pattern Recognition Preprocessors', In *International Conference on Microwaves, Circuit Theory, and Information Theory*, 1964
43. L. Kaufman and P.J. Rousseeuw, *Finding Groups in Data: an Introduction to Clustering Analysis*, Academic Press, San Diego, CA, 1990
44. E. Rasmussen, *Clustering Algorithms*, In *Information Retrieval: Data Structures and Algorithms*, W.B. Frakes and R. Baeza-Yates, Eds., Prentice-Hall, Inc., Upper Saddle River, NJ, pp. 419-442, 1992
45. N. Jardine and C.J. Rijsbergen, 'The Use of Hierarchical Clustering in Information Retrieval', *Information Storage and Retrieval*, **Vol.7**, pp. 217-240, 1971
46. R.S. Bennett, 'The Intrinsic Dimensionality of Signal Collections', *IEEE Transactions on Information Theory*, **Vol.15**, pp. 517-525, 1966
47. T. Zhang and R. Ramakrishnan and M. Livny, 'BIRCH: An efficient data clustering method for very large databases', *SIGMOD Rec.*, **Vol.25 No.2**, pp. 103-114, 1996
48. M. F. Porter, 'An Algorithm for Suffix Stripping', *Program*, **Vol.14 No.3**, pp. 130-137, 1980