# Document Clustering with Omni-Directional Data Similarity Process

**K. Lakshmi**

Department of Computer Science, Idhaya College for Women, Kumbakonam, Tamil Nadu, India

*Corresponding Author: luxmi.com@gmail.com*

*Abstract*— Most of the clustering techniques must presume some cluster relationship relating to the data thing. Similarity among some items is usually defined clearly or sometimes absolutely. In this paper, is an introduction to some novel reference centered similarity gauge and two related clustering approaches. The significant difference between an old-fashioned dissimilarity/similarity gauge and the approach considered in this paper is how the former uses simple single standpoint. In the existing approach it considers the origin, while the latter utilizes a number of reference details, which are objects assumed not to ever be inside the same cluster while using two things being scored. Using several reference details, more useful assessment of similarity could be possibly achieved. In document clustering two qualification functions are proposed and is determined by the fresh measure. The above functions are being examined along with frequently used clustering based algorithms which use other well known similarity measures in various document collections in order to verify the approach under consideration in this paper.

*Keywords*— Document Clustering, Correlation Measure, Similarity Measure, Data Mining

## I. INTRODUCTION

Clustering is among the most useful and essential areas within data mining. The objective of the clustering is always to find implicit structures within the data and organize them into some important sub groups intended for further study and analysis. There are many clustering algorithms published each year. They are usually proposed intended for very distinct research areas, and formulated using completely different techniques and approaches. According to the recent study k- means method still remains among the top 10 information mining algorithms, these days even though it was introduced half a decade earlier. It could be the most commonly used partitioned clustering algorithm used [1][2]. Another recent scientific talk denotes that k-means would be one of the best algorithms that are generally used by the practitioners from the respected fields. Needless to cover, k-means has many basic cons, such as sensitiveness in initialization and sizing cluster and its performance is naturally bad compared to other state-of-the-art algorithms in many domains. Despite that, it's simplicity, capability and scalability shall be the reasons to its tremendous recognition. Usability and performance are the main parameters in choosing an algorithm in most of the application scenarios as compared to performance & complexity. K-means approach easily combines with other techniques in more substantial systems with reasonable outcomes.

A common approach to the clustering problem is usually to treat it just as one optimization procedure. A best partition is available by optimizing a unique function connected with similarity (or distance) amongst data. Basically, there is definitely an implicit assumption which the true implicit structure connected with data could be correctly described because of the similarity system defined in addition to embedded within the clustering qualifying measure function. This is the reason why, the effectiveness involving clustering algorithms under such approach which depends on the appropriateness of the similarity measure is based on type of data available. For instance, the initial k-means features sum-of-squared-error function uses Euclidean distance. In an exceptionally sparse as well as high dimensional area like text documents, spherical k implies, the cosine similarity is more suited as compared to Euclidean distance [3].

## II. PROPOSED SYSTEM

### A. SimilarityMeasure

The particular cosine similarity might be expressed in a form applying Sim($d_i$, $d_j$), in which vector represents the origin point. According to the current formulae, our measure requires reference stage. The similarity between a couple of documents $d_i$ in addition to $d_j$ is decided w. r. t. the angle between two items when origin is considered as reference point. To build a new type of similarity, it is possible to use lot more than just one point involving reference. We may have a additional accurate assessment of exactly how close or perhaps distant a couple of things are. The huge reference items are suggested accordingly because of the difference vectors ($d_i − d_h$) in addition to ($d_j − d_h$).

The similarity of a couple of documents di in addition to dj signify that the inside same cluster pertains to the common properties of similarities measured relatively from the references of other clusters outside[4][5].

### B. Multi-Reference pointSimilarity

We call this module this Multi-Reference stage based Likeness, or MVS. Out of this point onwards, we tend to denote the suggested similarity calculation between 2 document vectors di in addition to dj simply is by MVS (di, dj). The MVS form is dependent upon a particular formulation on the individual similarities from the sum. If the relative likeness is identified by dot-product on the difference vectors, we now have: The likeness between two points di in addition to dj interior cluster Sr, referring to a point dh which is outside the present cluster of di, dj is equal to the product of cosine angle between di,dj and also the Euclidean distances from dh to these two points. This definition will be based upon the assumption that dh just isn't in the same cluster together with di, dj. Small distances di−dh in addition to dj −dh tend to have higher the chance that dh is in fact in the same cluster together with di, dj. Therefore, via these distances, we also supplies a measure of inter cluster dissimilarity, since points di in addition to dj fit in with cluster Sr, whereas dh belongs to an external cluster. The overall similarity between di in and dj depends on taking average of overall reference points not belonging to cluster Sr. It is possible to argue that when large number of reference points are considered, there may be a chance of getting mistaken information just as it can happen with the single reference point. Nonetheless, given a reasonable amount of reference points and their particular variety, it is reasonable to assume that most of them are useful. For this reason, the effect of mistaken reference points are reduced and ARE constrained by the average move. It is obvious that this multi reference similarity offers more informative assessment of similarity than the single reference point dependent similarity[6].

We call this module the Multi-Reference point based Similarity, or MRP. From this point onwards, we will denote the proposed similarity measure between two document vectors *di* and *dj*by MRS (*di, dj*). The MVS is defined as
MRS($d_i, d_j$ |$d_i$, $d_j$ Ɛ $S_r$) = (1/n-nr) $\sum$ cos($d_i − d_h, d_j − d_h$) ||$d_i − d_h$|| ||$d_j − d_h$||

   Where $d_i, d_j$ are two points in the same cluster .
      n is the total no of document
      $n_r$ is the no of documents in cluster _r'
      S is the set of all documents
      $d_h$ is the point outside the cluster
 *Algorithm:* Multi-ReferenceSimilarity
      Step 1: Retrieve all Documents
      Step 2: Compute relativesimilarity
      Step 3: Check for similarity inside cluster Sr,
      Step 4: Compute Euclidean distances
      Step 5: Compute cluster size-weighted

   Step 6: Compute average pair wise similarities
    Step 7: Compute intra-cluster similarity measure
   Step 8: Compute inter-cluster similarity measure
   Step 9: Compute similarity between each
         document vector and centroid

### C. ValidityComputation

For every sort of likeness measure, a likeness matrix A is established. For CS, this is simple, as aij = dtidj. The process for developing MVS matrix is first, the particular outer blend

w. r. t. each class is established. Then, for every single row ai where, i = 1,..., n, if the set of documents di and dj, where t = 1,..., n will be in the same class thus, aij is calculated. Normally, dj is assumed to stay di's type, and aij is calculated. After matrix A is shaped, the procedure is employed to receive its validity report. For each and every document di matching to line $a_i$ of $a_n$, we pick out qr files closest to di. The worth of qr is chosen reasonably as percentage of the length of the type r that contains di, exactly where percentage ∈ [(0, 1]. Then, validity w.r.t. *di* is calculated by the fraction of these *qr*documents having the same class label with *di*, The final validity is determined by averaging over all the rows of *A*. It is clear that validity score is bounded within 0 and 1. The higher validity score a similarity measure has, the more suitable it should be for the clustering task.[7].

### D. Clustering Criteria's

Having defined our similarity criteria, we now formulate our clustering qualifying measure functions. The first function, called IR, this would be the cluster size-weighted amount of average pairwise parallels of documents from the same cluster. We would want to transform this objective functionality into several suitable form to get facilitate the optimization procedure for being simple and fast. When comparing F with all the min-max lower, both functions would secure the two terms i.e. intra-cluster similarity measure and inter-cluster similarity measure. On the other hand, while the intention of min-max cut is to minimize this inverse relation between the two of these terms, our aim is to maximize their weighted difference. This difference term is resolute for every single cluster. They're weighted by the inverse from the cluster's dimensions, before summed up overall groups. One issue is that this formulation is supposed to be very sensitive in order to cluster dimensions. It shows up that IR's performance dependency on the value connected with α. The qualifying measure function yields relatively excellent clustering effects for α ∈ (0, 1). Inside the formulation connected with IR, a cluster quality is actually measured by the average pair wise similarity between the vectors/points within the same cluster. On the other hand, such an approach can result in sensitiveness towards size and also tightness from the clusters. Using CS, for example, pairwise similarity of documents in a sparse

cluster is generally smaller as compared to those in a dense cluster. To reduce this, a different approach is to consider similarity between every single document vector and its particular cluster's centroid alternatively[8][9].

## III.    RESULTS

The concept of this paper is implemented and different outcomes are illustrated below, the projected work is built on Java platform on Pentium-IV PC with 60 GB hard-disk and 1GB RAM. The implementation reveals that the efficiency has improved on differentReutersDatasets.TheFig1,Fig2, Fig 3, Fig 4 and Fig 5 shows the evaluation of the results[10][11].



Figure 1 Comparison of 200 documents in clusters.

The above graph Fig. 1 represents the number of documents clustered, we have taken 250 documents from reuters dataset and clustered them into 5 clusters. Similarly in Fig. 2 we have taken 500 documents from reuters dataset and clustered them in 5 clusters It is observed hat the documents in clusters changes with respect to cosine similarity and multi reference similarity[12][13].
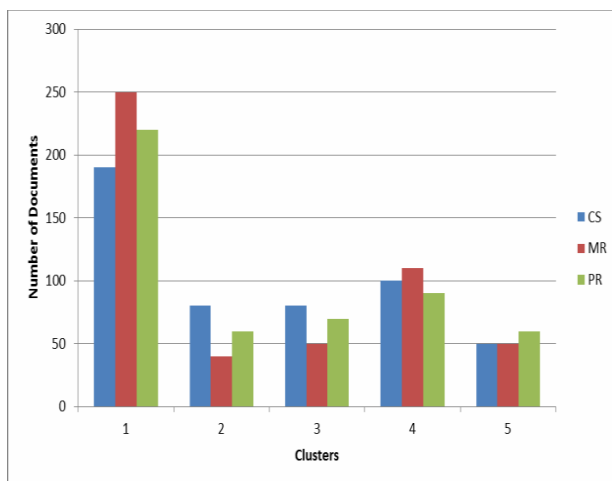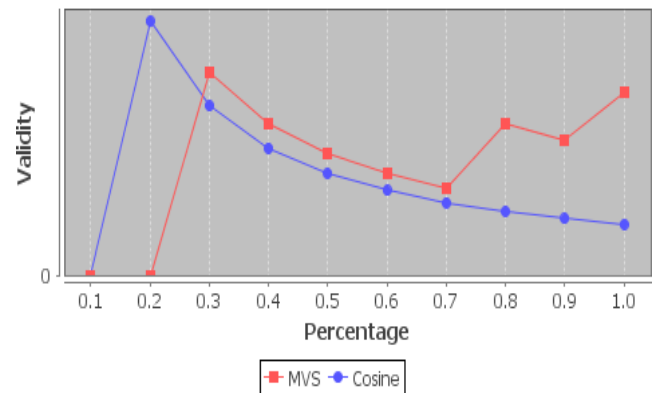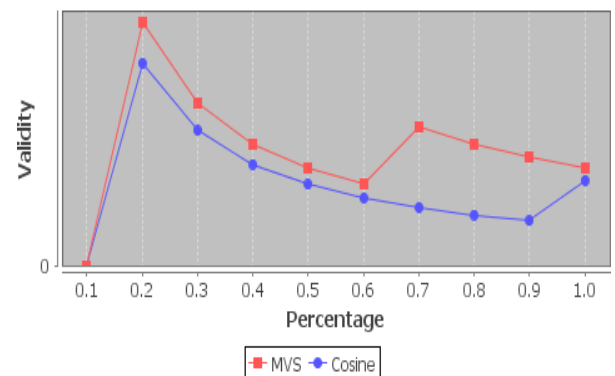


Figure 2 Comparison of 500 documents in clusters



Figure 3 Validity graph for 25documents

In the above graph Fig. 3, we compare the validity score for various percentages ranging from 0.1 to 1.0 in steps of 0.1. The above graph is plotted for 250documents.



Figure 4 Validity graph for 50documents

In the above graph Fig. 4, we compare the validity score for various percentages ranging from 0.1 to 1.0 in steps of 0.1. The above graph is plotted for 500documents.
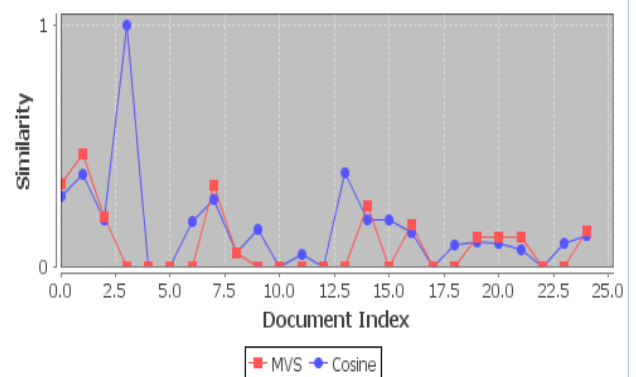


Figure 5 Comparison of Validity graph

In the above graph Fig. 5 depicts the comparision of similarity between cosine and multi reference for 250 documents.

## IV.    CONCLUSION AND FUTURE SCOPE

In this paper, we propose a Multi-Reference stage based Similarity measuring procedure, named MVS. Theoretical research and empirical cases depicts that MVS is potentially better for text message documents than the popular cosine likeness. Based on about MVS approach, a couple criterion characteristics, IR and IV, and their respective clustering algorithms, MVSC-IR and MVSC-IV, have been introduced. Compared the existing and suggested approaches using document datasets and under different evaluation metrics, the proposed algorithms show significantly enhanced clustering effectiveness.

The key contribution of the paper may be the fundamental concept of similarity measure from several reference factors. Future strategies could makes use of same theory, and might determine alternative forms for the relative likeness. This paper aims at partitioning clustering involving documents. Down the road, it would certainly also become possible to apply the proposed criterion capabilities for hierarchical clustering algorithms also. Lastly, we demonstrate the effective application of MVS as well as clustering algorithms pertaining to text files. It could well be interesting to explore in a direction that work on other forms of sparse in addition to high-dimensional files.

## REFERENCES

[1]   X. Wu, V. Kumar, J. Ross Quinlan, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu, P. S. Yu, Z.-H. Zhou, M.steinbach,D.J.Hand,andD.Steinberg,―Top10algorithmsindata mining,‖Knowl. Inf. Syst., vol. 14, no. 1, pp. 1–37, 2007.

[2]   HChimand X. Deng, ―Efficient phrase-based document similarity for clustering,‖ IEEE Trans. on Knowl. and Data Eng., vol. 20, no. 9, pp. 1217–1229,2008.

[3]   D. Lee andJ. Lee, ―Dynamic dissimilarity measure for support based clustering,‖ IEEE Trans. on Knowl. and Data Eng., vol. 22, no. 6, pp. 900–905,2010.

[4]   P. Lakkaraju, S. Gauch, andM. Speretta, ―Document similarity based on concept tree distance,‖ in Proc. of the 19th ACM conf. on Hypertext and hypermedia, 2008, pp.127–132.

[5]   D. Ienco, R. G. Pensa, andR. Meo, ―Context-based distance learning for categorical data clustering,‖ in Proc. of the 8th Int. Symp. IDA, 2009, pp.83–94.

[6]   I. Guyon, U.von Luxburg, and R. C. Williamson, ―Clustering: Science or Art?‖ NIPS'09 Workshop on Clustering Theory,2009.

[7]   E. Pekalska, A. Harol, R. P. W. Duin, B. Spillmann, and H. Bunke,―Non-Euclideanornon-metricmeasurescanbeinformative,‖ in Structural, Syntactic, and Statistical Pattern Recognition, ser. LNCS, vol. 4109, 2006, pp.871–880.

[8]   M.Pelillo, ―What is a cluster?Perspectivesfromgametheory,‖in Proc. of the NIPS Workshop on Clustering Theory,2009.

[9]   I. Dhillon and D. Modha, ―Concept decompositionsfor large sparse text data using clustering,‖ Mach. Learn., vol. 42, no. 1-2, pp. 143–175, Jan2001.

[10]  S. Zhong, ―Efficient online spherical K-means clustering,‖in IEEE IJCNN, 2005, pp. 3180–3185.