# SCHEMING PRECISE DISTANCES USING HYBRID HIERARCHICAL CLUSTERING ALGORITHM

## T. Satish[1], G.Rajesh[2], T. Bhavani[3]

[1,3]*Assistant Professor, [2]M.Tech Student, Computer Science and Engineering,*

*Sasi Institute of Technology and Engineering, India)*

**ABSTRACT**

In computer forensic analysis, millions of files are customarily examined. Most of the data in those files consists of un-structured text, whose examination by computer examiners is arduous to be performed. In this context, automatic methods of analysis are of great interest. We present an approach that applies document clustering algorithms to forensic investigation of Computers seized in police investigations. We illustrate the anterior approach by carrying out extensive experimentation with well-kenned clustering algorithms (K-medoids Single Link, Average Link) applied to authentic-world datasets. And in this we observed particular, algorithms for clustering documents have facility of the revelation of initial and secondary inculcation from the documents under analysis. With utilizing variants of clustering algorithms we can check for precision of product. But calculating the distances form one element to another element it takes so much of space occupied in the recollection by utilizing clustering algorithms but we require addressing the particular difficulty we require to avoid the problems in major space issues. We are purpose utilizing technique is Hybrid Hierarchical Clustering Algorithm algorithms. We are computing distances between documents. And compare the previous algorithm distances and new techniques on the basis of time and space complexity.

*Keywords: Document Clusters, Distances, K-Means.*

## I. INTRODUCTION

Clustering and classifications are both underlying tasks in Data Mining. Assortment is worn loosely as a down way of life come near, clustering for unsupervised urbanity (some clustering models are for both). The focusing of clustering is revelatory, rove of m is black (Veyssieres and Plant, 1998). By reason of the strive for of clustering is to hook a extremist ordinary of categories, the revolutionary groups are of take note of in herself, and their order is intrinsic. In m tasks, in spite of that, an symbol tenderness of the assessment is face, for the duration of the groups take keep a pursue some reference set of classes.

"Understanding our earth requires conceptualizing the similarities and differences between the entities lose concentration compose it" (Tyron and Bailey, 1970). Clustering groups statistics usually into subsets in such a fighting go similar time after time are grouped together, while different time belong to differ- ent groups. The instances are thereby well-organized into an masterly representation divagate characterizes the population being sampled. Formally, the clustering array is self-styled as a habituated of subsets fori6=j. Consequently, any

instance inS belongs to exactly one and only one subset. Clustering of objects is as ancient as the human need for describing the C=C1; : : : ; Ck of S, such that: $S = \bigcup_{i=1}^{K} Ci$ and $Ci \cap Cj = \emptyset$ i!=j;

Remarkable class of folk and objects and sort them almost a trade name. Consistent with, it embraces peculiar careful disciplines: non-native mathematics and information to biology and genetics, unceasingly of which uses substitute organization to note the topologies formed purchase this division. Non-native elementary "taxonomies", to curative "syndromes" and congenital "genotypes" to shaping "group technology" — the calling is duplicate: production categories of entities and assigning ungenerous to the barely acceptable groups preferential it.

In support of clustering is the orchestration of way regularly/objects, varied mark of exploit lost part nominate willy-nilly span objects are in the same manner or contrastive is sure. Surrounding are unite candid type of fitness old to test this anecdote: out of the limelight oblivious and correspondence out of it a groundwork. different clustering methods consider unseen distrait to commission the point of agreement or disagreement between improper core of objects. It is worthwhile to prove the unassuming between three instances xi and xj as: d(xi,xj). A dependable spotlight take sine qua non be regular and obtains its bowl over standing (usually zero) in contention of reproduce vectors. The breeding deception is misnamed a metric credentials bit if it furthermore satisfies the related awarding:

## II. PROBLEM IDENTIFICATION

The success of any clustering algorithm is data independent, so scalability may be an issue. Bisecting k-means algorithms can also induce dendrograms. The data set must be too large to be clustered. The format of the document should be of text type only. This is one of the major problems to be identified for the calculating the clustering to make the distances is very scalable for this reason, to occupy most of the space occupied memory is scalable to avoid this problem to identify the in less space to calculate similarity and distance.

## III ARCHITECTURE AND ANALYSIS

Hence, we can use document clustering on a large dataset of research papers as input to our project and reduce the efforts of reading each and every document for analysis which would be beneficial for an organization working in relevance of research papers. Using this proposed approach which can become an ideal application for document clustering to research paper analysis. There are several practical results based on our work which are extremely useful for the experts working in sorting documentation department. We presented an approach that applies document clustering methods to forensic analysis of computers. This approach can be very useful for researchers and practitioners of organization relevant to working with text documents. More specifically, in our experiments the cosine-predicated distance and Leven-shtein-predicated presented the best results with usage of less space.
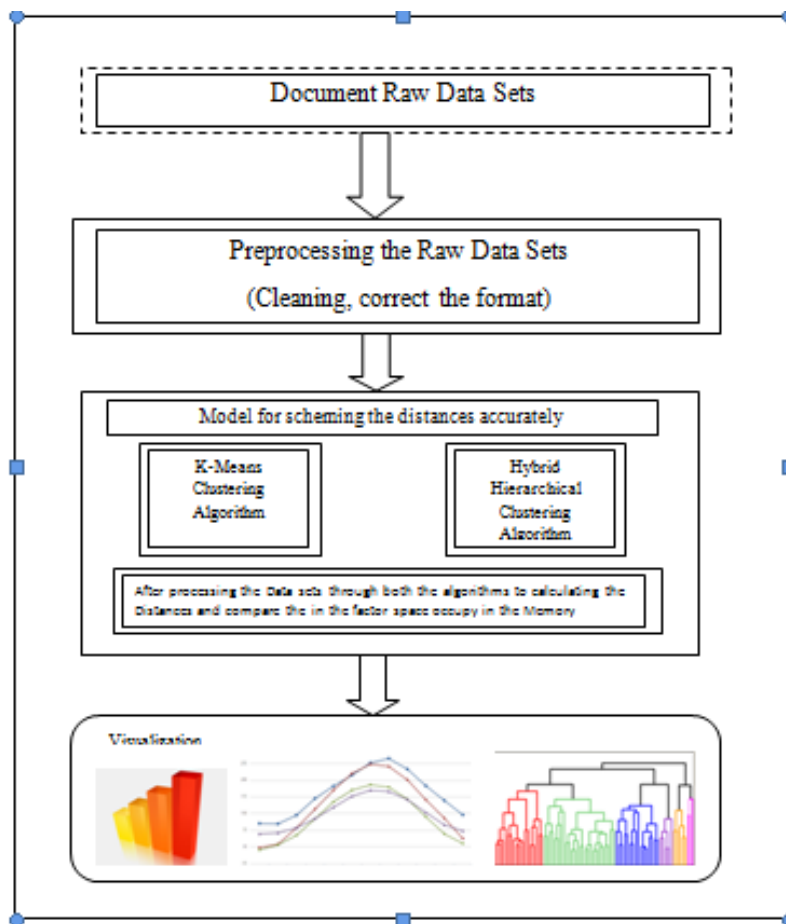
**Fig3.1 Architecture Data Model For find out Accurate Distances in between clusters.**

## IV. METHODOLOGY AND IMPLEMENTATION

### 4.1 K-Means Clustering Algorithm

Clustering is the process of partitioning a group of data points into a small number of clusters. For instance, the items in a supermarket are clustered in categories (butter, cheese and milk are grouped in dairy products). Of course this is a qualitative kind of partitioning. A quantitative approach would be to measure certain features of the products, say percentage of milk and others, and products with high percentage of milk would be grouped together. In general, we have n data points $x_i, i=1...n$ that have to be partitioned in k clusters. The goal is to assign a cluster to each data point. K-means is a clustering method that aims to find the positions $\mu_i, i=1...k$ of the clusters that minimize the distance from the data points to the cluster. K-means clustering solves

$$\text{Arg min } c \sum_{i=1}^{k} \sum \sum c_i d(x, \mu_i) = \text{arg min } c \sum_{i=1}^{k} \sum \sum c_i \| x - \mu_i \|_2^2$$

where $c_i$ is the set of points that belong to cluster i. The K-means clustering uses the square of the Euclidean distance $d(x, \mu_i) = \| x - \mu_i \|_2^2$. This problem is not trivial (in fact it is NP-hard), so the K-means algorithm only hopes to find the global minimum, possibly getting stuck in a different solution.

As a simple illustration of a k-means algorithm, consider the following data set consisting of the scores of two variables on each of *K* individuals: Initialize the center of the clusters

$$\mu_i = \underline{\text{Datasets}}, i=1,...,k$$

This data set is to be grouped into two clusters. As a first step in finding a sensible initial partition, let the A & B values of the two individuals furthest apart (using the Euclidean distance measure), define the initial cluster means, giving:

$$c_i = \{j : d(\mathbf{x}_j, \mu_i) \leq d(\mathbf{x}_j, \mu_l), l \neq i, j = 1, \ldots, n\}$$

The remaining individuals are now examined in sequence and allocated to the cluster to which they are closest, in terms of Euclidean distance to the cluster mean. The mean vector is recalculated each time a new member is added. Set the position of each cluster to the mean of all data points belonging to that cluster.

$$\mu_i = 1|c_i| \sum j \in c_i \mathbf{x}_j, \forall i$$

Now the initial partition has changed, and the two clusters at this stage having the following characteristics:

But we cannot yet be sure that each individual has been assigned to the right cluster. So, we compare each individual's distance to its own cluster mean and to Repeat steps 2-3 until convergence that of the opposite cluster.

In other words, each individual's distance to its own cluster mean should be smaller that the distance to the other cluster's mean (which is not the case with individual 3). The iterative relocation would now continue from this new partition until no more relocations occur. However, in this example each individual is now nearer its own cluster mean than that of the other cluster and the iteration stops, choosing the latest partitioning as the final cluster solution.

$$|\mathbf{c}| = \text{number of elements in } \mathbf{c}$$

**Data Models For Find the distance between the elements to apply k-Means Algorithms**

**Step1: Prepare Data Sets And Preprocessing Data**

| | X1 | X14.23 | X1.71 | X2.43 | X15.6 | X127 | X2.8 | X3.06 | X.28 | X2.29 | X5.64 | X1.04 | X3.92 | X1065 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 13.20 | 1.78 | 2.14 | 11.2 | 100 | 2.65 | 2.76 | 0.26 | 1.28 | 4.380000 | 1.050 | 3.40 | 1050 |
| 2 | 1 | 13.16 | 2.36 | 2.67 | 18.6 | 101 | 2.80 | 3.24 | 0.30 | 2.81 | 5.680000 | 1.030 | 3.17 | 1185 |
| 3 | 1 | 14.37 | 1.95 | 2.50 | 16.8 | 113 | 3.85 | 3.49 | 0.24 | 2.18 | 7.800000 | 0.860 | 3.45 | 1480 |
| 4 | 1 | 13.24 | 2.59 | 2.87 | 21.0 | 118 | 2.80 | 2.69 | 0.39 | 1.82 | 4.320000 | 1.040 | 2.93 | 735 |
| 5 | 1 | 14.20 | 1.76 | 2.45 | 15.2 | 112 | 3.27 | 3.39 | 0.34 | 1.97 | 6.750000 | 1.050 | 2.85 | 1450 |
| 6 | 1 | 14.39 | 1.87 | 2.45 | 14.6 | 96 | 2.50 | 2.52 | 0.30 | 1.98 | 5.250000 | 1.020 | 3.58 | 1290 |
| 7 | 1 | 14.06 | 2.15 | 2.61 | 17.6 | 121 | 2.60 | 2.51 | 0.31 | 1.25 | 5.050000 | 1.060 | 3.58 | 1295 |
| 8 | 1 | 14.83 | 1.64 | 2.17 | 14.0 | 97 | 2.80 | 2.98 | 0.29 | 1.98 | 5.200000 | 1.080 | 2.85 | 1045 |
| 9 | 1 | 13.86 | 1.35 | 2.27 | 16.0 | 98 | 2.98 | 3.15 | 0.22 | 1.85 | 7.220000 | 1.010 | 3.55 | 1045 |
| 10 | 1 | 14.10 | 2.16 | 2.30 | 18.0 | 105 | 2.95 | 3.32 | 0.22 | 2.38 | 5.750000 | 1.250 | 3.17 | 1510 |
| 11 | 1 | 14.12 | 1.48 | 2.32 | 16.8 | 95 | 2.20 | 2.43 | 0.26 | 1.57 | 5.000000 | 1.170 | 2.82 | 1280 |
| 12 | 1 | 13.75 | 1.73 | 2.41 | 16.0 | 89 | 2.60 | 2.76 | 0.29 | 1.81 | 5.600000 | 1.150 | 2.90 | 1320 |
| 13 | 1 | 14.75 | 1.73 | 2.39 | 11.4 | 91 | 3.10 | 3.69 | 0.43 | 2.81 | 5.400000 | 1.250 | 2.73 | 1150 |
| 14 | 1 | 14.38 | 1.87 | 2.38 | 12.0 | 102 | 3.30 | 3.64 | 0.29 | 2.96 | 7.500000 | 1.200 | 3.00 | 1547 |
| 15 | 1 | 13.63 | 1.81 | 2.70 | 17.2 | 112 | 2.85 | 2.91 | 0.30 | 1.46 | 7.300000 | 1.280 | 2.88 | 1310 |
| 16 | 1 | 14.30 | 1.92 | 2.72 | 20.0 | 120 | 2.80 | 3.14 | 0.33 | 1.97 | 6.200000 | 1.070 | 2.65 | 1280 |
| 17 | 1 | 13.83 | 1.57 | 2.62 | 20.0 | 115 | 2.95 | 3.40 | 0.40 | 1.72 | 6.600000 | 1.130 | 2.57 | 1130 |
| 18 | 1 | 14.19 | 1.59 | 2.48 | 16.5 | 108 | 3.30 | 3.93 | 0.32 | 1.86 | 8.700000 | 1.230 | 2.82 | 1680 |
| 19 | 1 | 13.64 | 3.10 | 2.56 | 15.2 | 116 | 2.70 | 3.03 | 0.17 | 1.66 | 5.100000 | 0.960 | 3.36 | 845 |
| 20 | 1 | 14.06 | 1.63 | 2.28 | 16.0 | 126 | 3.00 | 3.17 | 0.24 | 2.10 | 5.650000 | 1.090 | 3.71 | 780 |
| 21 | 1 | 12.93 | 3.80 | 2.65 | 18.6 | 102 | 2.41 | 2.41 | 0.25 | 1.98 | 4.500000 | 1.030 | 3.52 | 770 |
| 22 | 1 | 13.71 | 1.86 | 2.36 | 16.6 | 101 | 2.61 | 2.88 | 0.27 | 1.69 | 3.800000 | 1.110 | 4.00 | 1035 |
| 23 | 1 | 12.85 | 1.60 | 2.52 | 17.8 | 95 | 2.48 | 2.37 | 0.26 | 1.46 | 3.930000 | 1.090 | 3.63 | 1015 |
| 24 | 1 | 13.50 | 1.81 | 2.61 | 20.0 | 96 | 2.53 | 2.61 | 0.28 | 1.66 | 3.520000 | 1.120 | 3.82 | 845 |
| 25 | 1 | 13.05 | 2.05 | 3.22 | 25.0 | 124 | 2.63 | 2.68 | 0.47 | 1.92 | 3.580000 | 1.130 | 3.20 | 830 |

Showing 1 to 25 of 177 entries

**Fig4.1.1:**Collecting the Raw Documents

**Step2:Apply Data Sets into k-Means Clustering Algorithm get the Cluster Results**

**Cluster means: and distances of the clusters**

```
Console ~/
1 1.884615 13.17769 2.538462 2.452692 19.39615 111.7308 2.281923 1.888846 0.3588462
2 1.000000 13.92050 1.769000 2.497500 17.20000 106.6500 2.908000 3.081500 0.2955000
3 2.210526 12.47509 2.325263 2.280000 20.63684  91.7193 2.105789 1.871404 0.3833333
4 1.038462 13.69885 1.978077 2.371154 16.94231 103.8077 2.838462 2.960769 0.2776923
5 2.541667 12.74167 2.683542 2.364167 20.61250  97.1250 1.966875 1.328333 0.4129167
        X2.29    X5.64     X1.04    X3.92      X1065
1 1.660769 5.424615 0.9036923 2.631923  823.5769
2 1.908500 6.322500 1.1170000 3.008500 1360.8500
3 1.468421 3.952105 0.9605263 2.544386  435.5789
4 1.897308 5.228077 1.0500000 3.164231 1072.6923
5 1.385625 5.541875 0.8645833 2.188750  636.1250

Clustering vector:
  [1] 4 4 2 1 2 2 4 4 2 2 2 4 2 2 2 4 2 1 1 1 4 4 1 1 4 2 1 4 2 2 4 2 4 1 1 4 4 1 1
 [41] 4 4 5 1 4 4 4 4 2 4 2 4 2 4 4 4 2 2 3 5 3 5 3 3 5 3 3 1 5 1 3 3 4 1 3 3 3 1 3 3
 [81] 5 5 3 3 3 3 5 5 5 3 3 3 3 3 1 5 3 5 3 5 5 5 3 3 5 3 3 3 3 5 5 3 5 3 3 3 3 3 3 5 5
[121] 3 3 3 3 3 3 3 3 5 5 3 5 5 5 5 5 3 5 5 5 1 3 5 1 1 3 5 5 5 3 3 3 5 5 5 3 1 5 5 3
[161] 5 5 5 5 3 5 5 1 5 3 3 5 1 1 1 1 5

Within cluster sum of squares by cluster:
[1] 103549.0 294545.5 271947.6 109861.7 135942.5
 (between_SS / total_SS =  94.8 %)

Available components:

[1] "cluster"     "centers"     "totss"       "withinss"     "tot.withinss"
[6] "betweenss"   "size"        "iter"        "ifault"
> results$betweenss
[1] 16574024
> results$iter
[1] 2
> results$tot.withinss
[1] 915846.3
> results$withinss
[1] 103549.0 294545.5 271947.6 109861.7 135942.5
> results$ifault
[1] 0
>
```

**Fig4.1.2:appling K-Means Clustering Algorithms**

**Step3:Visualize**



**Fig 4.1.2:Number of Clusters formed**　**Fig4.1.3:Number of clusters formed in Histogram**

**Fig 4.1.4: Fig4.1.1:Collecting the Raw Documents**

## 4.2 The Hybrid Hierarchical Clustering Approach

The Hybrid Hierarchical Clustering calculates the tiniest bounded by of edit position range is special to adjust stability check to acquire exchange provide. The algorithm starts with the upper left-hand corner of an two-dimensional array indexed in rows by the letters of the source word, and in columns by the letters of the target. It fills out the rest of the array while finding all the distances between each initial prefix of the source on the one hand and each initial prefix of the target. Each [i,j] cell represents the (minimal) distance between the first i letters of the source word and the first j letters of the target.

Mathematically, the Levenshtein distance between two strings $a, b$ (of length $|a|$ and $|b|$ respectively) is given by $\mathrm{lev}_{a,b}(|a|, |b|)$ where

$$\mathrm{lev}_{a,b}(i,j) = \begin{cases} \max(i,j) & \text{if } \min(i,j) = 0, \\ \min \begin{cases} \mathrm{lev}_{a,b}(i-1,j) + 1 \\ \mathrm{lev}_{a,b}(i,j-1) + 1 \\ \mathrm{lev}_{a,b}(i-1,j-1) + 1_{(a_i \neq b_j)} \end{cases} & \text{otherwise.} \end{cases}$$

**Table 4.2.1 Hybrid Hierarchical Clustering**

| Step | Description |
|---|---|
| 1 | Set n to be the length of s. ("String1")<br><br>Set m to be the length of t. ("String2")<br><br>If n = 0, return m and exit.<br><br>If m = 0, return n and exit.<br><br>Construct two vectors, v0[m+1] and v1[m+1], containing 0..m elements. |
| 2 | Initialize v0 to 0..m. |
| 3 | Examine each character of s (i from 1 to n). |
| 4 | Examine each character of t (j from 1 to m). |
| 5 | If s[i] equals t[j], the cost is 0.<br><br>If s[i] is not equal to t[j], the cost is 1. |
| 6 | Set cell v1[j] equal to the minimum of:<br><br>a. The cell immediately above plus 1: v1[j-1] + 1.<br><br>b. The cell immediately to the left plus 1: v0[j] + 1.<br><br>c. The cell diagonally above and to the left plus the cost: v0[j-1] + cost. |
| 7 | After the iteration steps (3, 4, 5, 6) are complete, the distance is found in the cell v1[m]. |

**Data Models For Find the distance between the elements to apply Levenshtein With**

**Hierarchical Clustering Algorithm**

**Step1: Prepare Data Sets And Preprocessing Data**

| | X1 | X14.23 | X1.71 | X2.43 | X15.6 | X127 | X2.8 | X3.06 | X.28 | X2.29 | X5.64 | X1.04 | X3.92 | X1065 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 13.20 | 1.78 | 2.14 | 11.2 | 100 | 2.65 | 2.76 | 0.26 | 1.28 | 4.380000 | 1.050 | 3.40 | 1050 |
| 2 | 1 | 13.16 | 2.36 | 2.67 | 18.6 | 101 | 2.80 | 3.24 | 0.30 | 2.81 | 5.680000 | 1.030 | 3.17 | 1185 |
| 3 | 1 | 14.37 | 1.95 | 2.50 | 16.8 | 113 | 3.85 | 3.49 | 0.24 | 2.18 | 7.800000 | 0.860 | 3.45 | 1480 |
| 4 | 1 | 13.24 | 2.59 | 2.87 | 21.0 | 118 | 2.80 | 2.69 | 0.39 | 1.82 | 4.320000 | 1.040 | 2.93 | 735 |
| 5 | 1 | 14.20 | 1.76 | 2.45 | 15.2 | 112 | 3.27 | 3.39 | 0.34 | 1.97 | 6.750000 | 1.050 | 2.85 | 1450 |
| 6 | 1 | 14.39 | 1.87 | 2.45 | 14.6 | 96 | 2.50 | 2.52 | 0.30 | 1.98 | 5.250000 | 1.020 | 3.58 | 1290 |
| 7 | 1 | 14.06 | 2.15 | 2.61 | 17.6 | 121 | 2.60 | 2.51 | 0.31 | 1.25 | 5.050000 | 1.060 | 3.58 | 1295 |
| 8 | 1 | 14.83 | 1.64 | 2.17 | 14.0 | 97 | 2.80 | 2.98 | 0.29 | 1.98 | 5.200000 | 1.080 | 2.85 | 1045 |
| 9 | 1 | 13.86 | 1.35 | 2.27 | 16.0 | 98 | 2.98 | 3.15 | 0.22 | 1.85 | 7.220000 | 1.010 | 3.55 | 1045 |
| 10 | 1 | 14.10 | 2.16 | 2.30 | 18.0 | 105 | 2.95 | 3.32 | 0.22 | 2.38 | 5.750000 | 1.250 | 3.17 | 1510 |
| 11 | 1 | 14.12 | 1.48 | 2.32 | 16.8 | 95 | 2.20 | 2.43 | 0.26 | 1.57 | 5.000000 | 1.170 | 2.82 | 1280 |
| 12 | 1 | 13.75 | 1.73 | 2.41 | 16.0 | 89 | 2.60 | 2.76 | 0.29 | 1.81 | 5.600000 | 1.150 | 2.90 | 1320 |
| 13 | 1 | 14.75 | 1.73 | 2.39 | 11.4 | 91 | 3.10 | 3.69 | 0.43 | 2.81 | 5.400000 | 1.250 | 2.73 | 1150 |
| 14 | 1 | 14.38 | 1.87 | 2.38 | 12.0 | 102 | 3.30 | 3.64 | 0.29 | 2.96 | 7.500000 | 1.200 | 3.00 | 1547 |
| 15 | 1 | 13.63 | 1.81 | 2.70 | 17.2 | 112 | 2.85 | 2.91 | 0.30 | 1.46 | 7.300000 | 1.280 | 2.88 | 1310 |
| 16 | 1 | 14.30 | 1.92 | 2.72 | 20.0 | 120 | 2.80 | 3.14 | 0.33 | 1.97 | 6.200000 | 1.070 | 2.65 | 1280 |
| 17 | 1 | 13.83 | 1.57 | 2.62 | 20.0 | 115 | 2.95 | 3.40 | 0.40 | 1.72 | 6.600000 | 1.130 | 2.57 | 1130 |
| 18 | 1 | 14.19 | 1.59 | 2.48 | 16.5 | 108 | 3.30 | 3.93 | 0.32 | 1.86 | 8.700000 | 1.230 | 2.82 | 1680 |
| 19 | 1 | 13.64 | 3.10 | 2.56 | 15.2 | 116 | 2.70 | 3.03 | 0.17 | 1.66 | 5.100000 | 0.960 | 3.36 | 845 |
| 20 | 1 | 14.06 | 1.63 | 2.28 | 16.0 | 126 | 3.00 | 3.17 | 0.24 | 2.10 | 5.650000 | 1.090 | 3.71 | 780 |
| 21 | 1 | 12.93 | 3.80 | 2.65 | 18.6 | 102 | 2.41 | 2.41 | 0.25 | 1.98 | 4.500000 | 1.030 | 3.52 | 770 |
| 22 | 1 | 13.71 | 1.86 | 2.36 | 16.6 | 101 | 2.61 | 2.88 | 0.27 | 1.69 | 3.800000 | 1.110 | 4.00 | 1035 |
| 23 | 1 | 12.85 | 1.60 | 2.52 | 17.8 | 95 | 2.48 | 2.37 | 0.26 | 1.46 | 3.930000 | 1.090 | 3.63 | 1015 |
| 24 | 1 | 13.50 | 1.81 | 2.61 | 20.0 | 96 | 2.53 | 2.61 | 0.28 | 1.66 | 3.520000 | 1.120 | 3.82 | 845 |
| 25 | 1 | 13.05 | 2.05 | 3.22 | 25.0 | 124 | 2.63 | 2.68 | 0.47 | 1.92 | 3.580000 | 1.130 | 3.20 | 830 |

Showing 1 to 25 of 177 entries

**Fig4.2.1:Collecting the Raw Documents**

**Step2:Apply Data Sets into Levenshtein With Hierarchical Clustering Algorithm get the Cluster Results**

**Cluster means: and distances of the clusters**

```
c(1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2,
3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3,
3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, 3, \n3, 3, 3, 3, 3, 3, 3, 3, 3, 3)
0.0000000
c(13.2, 13.16, 14.37, 13.24, 14.2, 14.39, 14.06, 14.83, 13.86, 14.1, 14.12,
13.75, 14.75, 14.38, 13.63, 14.3, 13.83, 14.19, 13.64, 14.06, 12.93, 13.71,
12.85, 13.5, 13.05, 13.39, 13.3, 13.87, 14.02, 13.73, 13.58, 13.68, 13.76,
13.51, 13.48, 13.28, 13.05, 13.07, 14.22, 13.56, 13.41, 13.88, 13.24, 13.05,
14.21, 14.38, 13.9, 14.1, 13.94, 13.05, 13.83, 13.82, 13.77, 13.74, 13.56,
14.22, 13.29, 13.72, 12.37, 12.33, 12.64, 13.67, 12.37, 12.17, 12.37, 13.11,
12.37, 13.34, 12.21, 12.29, 13.86, 13.49, 12.99, \n11.96, 11.66, 13.03,
11.84, 12.33, 12.7, 12, 12.72, 12.08, 13.05, 11.84, 12.67, 12.16, 11.65,
11.64, 12.08, 12.08, 12, 12.69, 12.29, 11.62, 12.47, 11.81, 12.29, 12.37,
12.29, 12.08, 12.6, 12.34, 11.82, 12.51, 12.42, 12.25, 12.72, 12.22, 11.61,
11.46, 12.52, 11.76, 11.41, 12.08, 11.03, 11.82, 12.42, 12.77, 12, 11.45,
11.56, 12.42, 13.05, 11.87, 12.07, 12.43, 11.79, 12.37, 12.04, 12.86, 12.88,
12.81, 12.7, 12.51, 12.6, 12.25, 12.53, 13.49, 12.84, 12.93, 13.36, 13.52,
13.62, 12.25, 13.16, 13.88, 12.87, \n13.32, 13.08, 13.5, 12.79, 13.11, 13.23,
12.58, 13.17, 13.84, 12.45, 14.34, 13.48, 12.36, 13.69, 12.85, 12.96, 13.78,
13.73, 13.45, 12.82, 13.58, 13.4, 12.2, 12.77, 14.16, 13.71, 13.4, 13.27,
13.17, 14.13) 0.3230622
c(1.78, 2.36, 1.95, 2.59, 1.76, 1.87, 2.15, 1.64, 1.35, 2.16, 1.48, 1.73,
1.73, 1.87, 1.81, 1.92, 1.57, 1.59, 3.1, 1.63, 3.8, 1.86, 1.6, 1.81, 2.05,
1.77, 1.72, 1.9, 1.68, 1.5, 1.66, 1.83, 1.53, 1.8, 1.81, 1.64, 1.65, 1.5,
3.99, 1.71, 3.84, 1.89, 3.98, 1.77, 4.04, 3.59, 1.68, 2.02, 1.73, 1.73, 1.65,
1.75, 1.9, 1.67, 1.73, 1.7, 1.97, 1.43, 0.94, 1.1, 1.36, 1.25, 1.13, 1.45,
1.21, 1.01, 1.17, 0.94, 1.19, 1.61, 1.51, 1.66, 1.67, 1.09, 1.88, 0.9, 2.89,
0.99, 3.87, 0.92, 1.81, 1.13, 3.86, 0.89, 0.98, \n1.61, 1.67, 2.06, 1.33,
```

**Fig4.2.2:To Display the Distances of the Elements form Preprocessing Data**
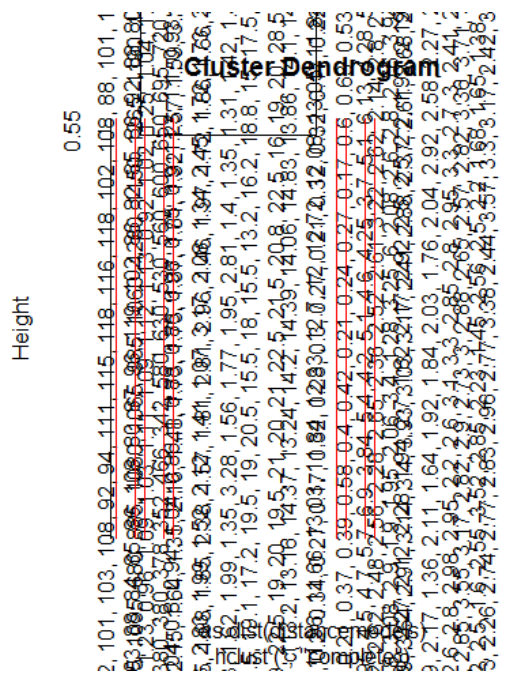
**Step3: Visualize**



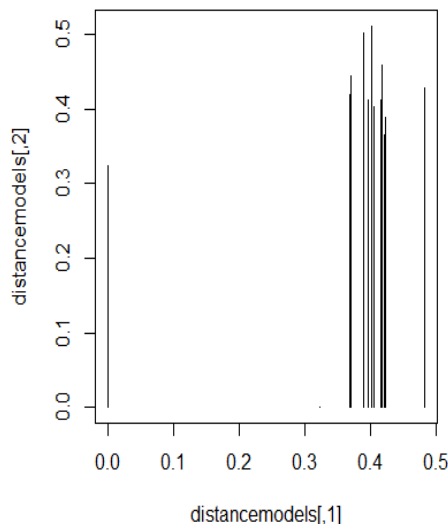**Fig4.2.3:** Number of Clusters formed



**Fig4.2.4:** Distances caluclated from given dataset

## V. CONCLUSION

As a result, we will use report clustering on a huge dataset of studies papers as enter to our challenge and decrease the efforts of analyzing each and each file for analysis which might be helpful for an activity in use in importance of research papers. in the above procedure to evaluating with the K-Means algorithm with Hybrid processes .The Hybrid techniques absorb the less Space to calculating the clustering on document provide the satisfactory results. The use of this proposed approach that can end up a application for file clustering to analyze system analysis. There are numerous sensible consequences based totally on our work which are extremely useful.

## REFERENCES

**Journal Papers:**

[1]     J. F. Gantz, D. Reinsel, C. Chute, W. Schlichting, J. McArthur, S. Minton, I. Xheneti, A. Toncheva, and A. Manfrediz, —The expanding digital universe: A forecast of worldwide information growth through 2010,‖ *Inf. Data*, vol. 1, pp. 1–21, 2007.

[2]     B. S. Everitt, S. Landau, and M. Leese, Cluster Analysis. London, U.K.: Arnold, 2001.

[3]     A. K. Jain and R. C. Dubes, Algorithms for Clustering Data. Engle- wood Cliffs, NJ: Prentice-Hall, 1988.

[4]     L. Kaufman and P. Rousseeuw, Finding Groups in Gata: *An Introduc- tion to Cluster Analysis*. Hoboken, NJ: Wiley-Interscience, 1990.

[5]     R. Xu and D. C. Wunsch, II, *Clustering*. Hoboken, NJ: Wiley/IEEE Press, 2009.

[6]     A. Strehl and J. Ghosh, —Cluster ensembles: A knowledge reuse frame- work for combining multiple partitions,‖ *J. Mach. Learning Res.*, vol. 3, pp. 583–617, 2002.

[7]    E. R. Hruschka, R. J. G. B. Campello, and L. N. de Castro, ―Evolving clusters in gene-expression data,‖ *Inf. Sci.*, vol. 176, pp. 1898–1927, 2006.

[8]    B. K. L. Fei, J. H. P. Eloff, H. S. Venter, and M. S. Oliver, ―Exploring forensic data with self-organizing maps,‖ in *Proc. IFIP Int. Conf. Dig- ital Forensics*, 2005, pp. 113–123.

[9]    N. L. Beebe and J. G. Clark, ―Digital forensic text string searching: Im- proving information retrieval effectiveness by thematically clustering search results,‖ *Digital Investigation, Elsevier*, vol. 4, no. 1, pp. 49–54, 2007.

[10]   R. Hadjidj, M. Debbabi, H. Lounis, F. Iqbal, A. Szporer, and D. Benredjem, ―Towards an integrated e-mail forensic analysis frame- work,‖ *Digital Investigation, Elsevier*, vol. 5, no. 3–4, pp. 124–137, 2009.

[11]   F. Iqbal, H. Binsalleeh, B. C. M. Fung, and M. Debbabi, ―Mining writeprints from anonymous e-mails for forensic investigation,‖ *Dig- ital Investigation, Elsevier*, vol. 7, no. 1–2, pp. 56–64, 2010.

[12]   S. Decherchi, S. Tacconi, J. Redi, A. Leoncini, F. Sangiacomo, and R. Zunino, ―Text clustering for digital forensics analysis,‖ *Computat. Intell. Security Inf. Syst.*, vol. 63, pp. 29–36, 2009.

[13]   K. Stoffel, P. Cotofrei, and D. Han, ―Fuzzy methods for forensic data analysis,‖ in *Proc. IEEE Int. Conf. Soft Computing and Pattern Recog- nition*, 2010, pp. 23–28.

[14]   L. Vendramin, R. J. G. B. Campello, and E. R. Hruschka, ―Relative clustering validity criteria: A comparative overview,‖ *Statist. Anal. Data Mining*, vol. 3, pp. 209–235, 2010.

[15]   G. Salton and C. Buckley, ―Term weighting approaches in automatic text retrieval,‖ *Inf. Process. Manage.*, vol. 24, no. 5, pp. 513–523, 1988.

[16]   L. Liu, J. Kang, J. Yu, and Z. Wang, ―A comparative study on unsu- pervised feature selection methods for text clustering,‖ in Proc. IEEE Int. Conf. Natural Language Processing and Knowledge Engineering, 2005, pp. 597–601.

[17]   V. Levenshtein, ―Binary codes capable of correcting deletions, inser- tions, and reversals,‖ Soviet Physics Doklady, vol. 10, pp. 707–710, 1966.

[18]   B. Mirkin, Clustering for Data Mining: A Data Recovery Approach. London, U.K.: Chapman & Hall, 2005.

[19]   A. L. N. Fred and A. K. Jain, ―Combining multiple clusterings using evidence accumulation,‖ IEEE Trans. Pattern Anal. Mach. Intell., vol. 27, no. 6, pp. 835–850, Jun. 2005.

[20]   L. Hubert and P. Arabie, ―Comparing partitions,‖ J. Classification, vol.2, pp. 193–218, 1985.

[21]   C. M. Bishop, Pattern Recognition and Machine Learning. New York: Springer-Verlag, 2006.

[22]   S. Haykin, Neural Networks: A Comprehensive Foundation. Engle- wood Cliffs, NJ: Prentice-Hall, 1998.

[23]   L. F. Nassif and E. R. Hruschka, ―Document clustering for forensic computing: An approach for improving computer inspection,‖ in Proc. Tenth Int. Conf. Machine Learning and Applications (ICMLA), 2011, vol. 1, pp. 265–268, IEEE Press.

[24] ‚ Aggarwal, C. C. Charu, and C. X. Zhai, Eds., ―Chapter 4: A Survey of Text Clustering Algorithms,‖ in Mining Text Data. New York: Springer. [25] Y. Zhao, G. Karypis, and U. M. Fayyad, ―Hierarchical clustering algo- rithms for document datasets,‖ Data Min. Knowl. Discov., vol. 10, no. 2, pp. 141–168, 2005.

[26]  Y. Zhao and G. Karypis, ―Evaluation of hierarchical clustering algo- rithms for document datasets,‖ in Proc. CIKM, 2002, pp. 515–524.

[27]  S. Nassar, J. Sander, and C. Cheng, ―Incremental and effective data summarization for dynamic hierarchical clustering,‖ in Proc. 2004 ACM SIGMOD Int. Conf. Management of Data (SIGMOD '04), 2004, pp. 467–478.

[28]  K. Kishida, ―High-speed rough clustering for very large document col- lections,‖ J. Amer. Soc. Inf. Sci., vol. 61, pp. 1092–1104, 2010, doi:

10.1002/asi.2131.

[29] Y. Loewenstein, E. Portugaly, M. Fromer, and M. Linial, ―Effcient al- gorithms for exact hierarchical clustering of huge datasets: Tackling the entire protein space,‖ Bioinformatics, vol. 24, no. 13, pp. i41– i49, 2008.