

High Dimensional Object Analysis Using Rough-Set Theory and Grey Relational Clustering Algorithm

Prashant Verma¹, Yogendra Kumar Jain²

Research Scholar, Computer Science and Engineering, Samrat Ashok Technological Institute, Vidisha (MP), India¹

Head of the Department, Computer Science and Engg., Samrat Ashok Technological Institute, Vidisha (MP), India²

Abstract: High dimensional feature selection and data assignment is an important feature for high dimensional object analysis. In this work, we propose a new hybrid approach of combining attribute reduction of the Rough-set theory with Grey relation clustering. Designing clustering becomes increasingly tougher task as the dimensionality of the data set increases. Previously constraint based clustering algorithms that satisfy user specified constraints have been used for high dimensional data sets. Such algorithms suffer from serious limitations and can introduce biases of the user, thus obscuring discovery of clusters and hidden relations in the data set. In this work, we transform the high relevance values into the same class using Grey relation to give an appropriate cluster of information, which we process through Rough set to reduce attributes. We use this approach to analyze the data of plant diversity from North America and find that ground elevation and species numbers can capture the most important attributes of the data set. This analysis of ecological data presents a proof of principal for the novel hybrid approach using Grey relational clustering and Rough set theory.

Keywords: RST, GRA, Rule Generation, High Dimensionality, Indiscernibility (IND).

1. INTRODUCTION

Clustering provide a better understanding of the data by dividing data point into clusters such that objects in the same cluster are similar[1], whereas objects in different cluster are dissimilar with respect to a given similarity measure[2]. Clustering of many algorithms has been studied for decades but in the age of data deluge conventional clustering algorithm is showing cracks and novel algorithms are needed. In the case of high dimensional data, a problem of clustering of data points that do not have enough feature relevance becomes a big problem[3]. Thus, data clustering in the case of high dimensional data poses two separate problems: (1) the search for relevant sub spaces[4] and (2) the detection of the final clusters

High dimensional data clustering algorithms could be categorized by their ways of dealing with local feature relevance. Subspace clustering algorithms employ dimension selection methods to form a subspace for each cluster. Subspace clustering algorithms can be both hard and soft. In hard clustering, where one datum point can belong to one and only one cluster, the performance subspace clustering algorithms is frequently hindered by the tough choice of relevant dimensions of clusters. Errors of missing relevant dimensions and inclusion of irrelevant dimensions also cause problems in hard subspace clustering. In hard subspace clustering algorithms, the selected dimensions of each cluster are viewed as equally important. However, in reality the dimensions of each subspace are usually not uniformly important in the same way for all the clusters.

In soft subspace clustering a datum point can belong to more than one cluster. While soft subspace clustering algorithms can remove irrelevant dimensions by not assigning a specific subspace for each cluster but it fails to deal with the problem of feature relevance[5]. The irrelevant dimensions, which are usually low weighted, tend to add noise to the procedures of finding cluster in these algorithms, leading to poor clustering results[6]. It seems that these kinds of algorithms could be adapted to include a dimension selection function by assigning some dimensions with 0 weights; nevertheless it's hard to determine which dimensions should be 0 weighted and until now there is no such a scheme. Moreover, there are usually a small number of relevant dimensions and very large number of irrelevant dimensions for each cluster.

Thus such a scheme is inefficient. However, if we perform dimension selection firstly and then perform dimension weighting, the computation can be largely reduced. To detect the final cluster, most high dimensional data clustering algorithms adopt a centroid-based approach, where initial centroids are established, followed by assigning data points to the closest centroid. Updating the centroids and reassigning data point according to some optimization criterion refines the clusters. From the above discussion we conclude that dimension selection, dimension weighting and data assignment (initial and reassignment) are three essential tasks for high dimension data clustering. High dimensional data clustering is a challenging science. Each underlying task is hard to solve and to add to the woes, the three tasks of dimension

selection, dimension weighting and data assignment are circularly dependent on each other.

2. PROBLEM DEFINITION

2.1 Introduction

The common theme of these problems is that when the dimensionality increases, the volume of the space increases so fast that the available data become sparse. This sparsity is problematic for any method that requires statistical significance. In order to obtain a statistically sound and reliable result, the amount of data needed to support the result often grows exponentially with the dimensionality. This led to the phrase “curse of dimensionality” by Richard E. Bellman, when considering problems in dynamic optimization. For distance functions and nearest neighbor search, recent research shows that data sets that are sparse due to high dimensionality can still be processed, unless there are too many irrelevant dimensions, while relevant dimensions can make some problems such as clustering actually easier. Any low-dimensional data space can trivially be turned into a higher-dimensional space by adding redundant or randomized dimensions, and in turn many high-dimensional data sets can be reduced to lower-dimensional data without significant information loss of information.

2.2 Dimensionality Reduction

Dimensionality Reduction is a process of reducing attributes from the data set. Within a data set there are exists superfluous information that is non-essential and this superfluous information contributes to the data complexity, increasing the time needed for analysis. There are several methods of dimensionality reduction, with the common ones being:

- Independent Component Analysis
- Principal Component Analysis (PCA)
- Probabilistic PCA (PPCA)
- The Kernel Trick
- Kernel PCA
- Canonical Correlation Analysis
- Linear Discriminant Analysis

I am briefly describing PCA as an example, as it is used primarily for linear data sets and I am also employing linear data set for my analysis. The Principal Component Analysis (PCA) is one of the dimension reduction methods consisting of the transfer of data to a new orthogonal basis, whose axes are oriented in the directions of the maximum variance of the input data set. The variance is maximum along the first axis of the new basis, while the second axis maximizes variance, subject to the first axis orthogonally, and so forth, the last axis having the least variance of all possible ones. Such transformation permits information to be reduced by rejecting the coordinates that correspond to the directions with a minimum variance. If one of the base vectors needs to be rejected, that should preferably be the vector along which the input data set is less changeable. In most cases, PCA does not guarantee that the selected first principal components will be the most adequate for

classification. One of the possibilities for selecting discriminative features from principal components is to apply rough sets theory.

3. PROPOSED METHOD

Many real-world data sets consist of a very high dimensional feature space. Clustering real-world data sets is often hampered by the so-called curse of dimensionality. Most of the common algorithms fail to generate meaningful results for clustering because of the inherent sparsity of the data space. Usually, clusters cannot be found in the original feature space because several features may be irrelevant for clustering. However, clusters are usually embedded in the lower dimensional subspaces. In addition, different sets of features may be relevant for different sets of objects. Thus, objects can often be clustered differently in varying subspaces of the original feature space. In this thesis, I am going to cluster high dimensional data using hybrid approach. I briefly describe the two approaches that I am hybridizing.

3.1 Rough Set

Rough set theory is used in various research areas, such as soft computing, machine learning, decision-making, data mining and KDD (Knowledge Data Discovery) for data analysis. Rough Set Theory is very helpful in reduction of dimensionality from high dimensional data sets. Rough Set theory was introduced by Zdzislaw Pawlak in the early 1980s [7]. Rough set is a fastest growing mathematical tool, which deals with intelligence and espionage data and data mining. Figure 1 summarizes the scheme of rough set.

Let consider there are information set $S = \langle U, A \rangle$ where U represents the set of non-empty finite objects $\{U_{i=1}^n U_i\}$ and A represents set of non-empty finite attributes. Where $\{a\}$ is the value of attribute generally represents as $\{a: U \rightarrow V_a\}$ for every $\{a \in A\}$. In any information system the set of attributes are the collection of conditional attribute $\{C\}$ and decision attribute $\{D\}$ hence we can represents the $\{A = C \cup D\}$ in information set equation $DT = \langle U, C \cup D \rangle$ is known as decision table. A decision table can be classified as a supervised learning, where the outcome of any system are well known and this posterior knowledge is well distinguished in an attribute that is called “Decision Attribute”. If $p \subseteq A$ then p is an association equivalence relation and this relation can also specified as a Indiscernible Relation. Assume $\delta = (U, A)$ is an information system, then any $p \subseteq A$ associated with equivalence class can be represents as $IND_\delta(p)$.

Now we are exploring two important terms in rough set theory:

Approximation defines when p is an association relation with attribute set A $\{p \subseteq A\}$ and $\{X \subseteq U\}$ can be approximated using the information in p construction. Lower Approximation: A lower approximation R_* represents the values which are surely belongs in set.

$R_* = \{U_x \in U \{p(x) \subseteq X\}\}$

Upper Approximation: An Upper approximation R^* represents those values which are possibly belongs in set.

$$R^* = \{U_x \in U \{p(x): p(x) \cap X \neq \emptyset\}$$

A boundary region in a Rough set is describe as those objects that can neither ruled nor ruled out as a member of target set X. represents as $R^* - R_*$ if there is an empty region then it look likes $R^* = R_*$. this situation belongs to crisp set if it does not happen that mean it's in Rough set.

3.2 Dependency of Attribute

Dependency in an attribute of similarly and can be extracted from relational data set. If all the values of any attribute A1 are uniquely determined by attribute A2 then we can say that attribute A1 totally depends on attribute A2 and this expression represents as $(A2 \rightarrow A1)$. We can also measure the degree of dependency which deviates between (0, 1). It can be easily seen that if D depends totally on A2, then $I(A2) \subseteq I(A1)$.

That means that the partition generated by A2 is finer than the partition generated by A1.

Notice that Dependency discussed above corresponds to that considered in relational databases. If A1 depends on the degree of k, $0 \leq k \leq 1$, on C, then.

$$\delta(A2, A1) = \frac{|Pos_{A2}(A1)|}{|U|}$$

Where $Pos_{A2}(A1) = \cup A2(X), X \in U/I(A1)$

3.4 Rough Set Reduct

Reduct in a rough set theory applied when attribute reduction is needed, when the information set are having dispensable attributes that are increasing unwanted weight of the information. Reduct reduces the dispensable attribute without changing its original classification [8]. Thus the reduct is the minimal subset of attribute that enables the classification of the elements.

$$core(T) = \cap Reduct(T)$$

Where $core(T)$ is set of all indispensable attribute of T and $Reduct(T)$ is the set of all superfluous elements.

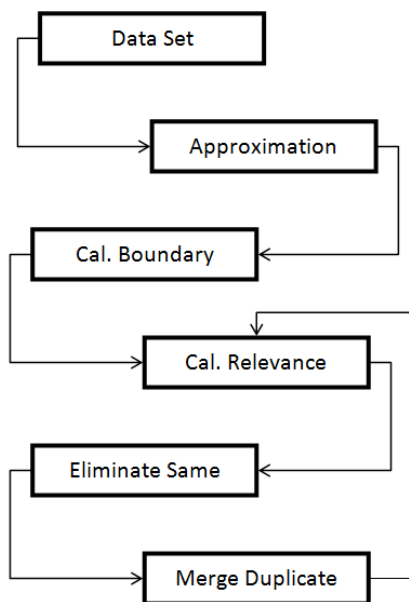


Figure 3.1: Process Diagram of Rough Set Dimensionality Reduction

ALGORITHM 1.0 –RST

```

    Begin
      Initialization : C=Conditional Attribute , D=
      Decision Attribute
      If(I(Q)==I(Q-{a})) Begin
        Then a= dispensable;
        Else a= Indispensable;
      End
      //Select Core
      Begin core(T) = ∩reduct(T)
      End
    End
  
```

3.5 Grey System:

Prof. Deng introduced the concept of Grey system theory in 1982. Grey system theory took the hypothetical black box and white box approach, representing unknown and known values respectively and introduced moderate values that are partially known and partially unknown as the grey system. In 1989, Prof. Deng proposed another theory on Grey Clustering Analysis (GCA), where he also described grey number, and grey equations. Grey relation analysis describes the relationship degree of objects, which extends the discrete sequence of values. Grey clustering relation explores the relation through hierarchal structure and has the flexibility in nature of classification, while exhibiting an effective performance[9].

3.6 Grey Relation Analysis

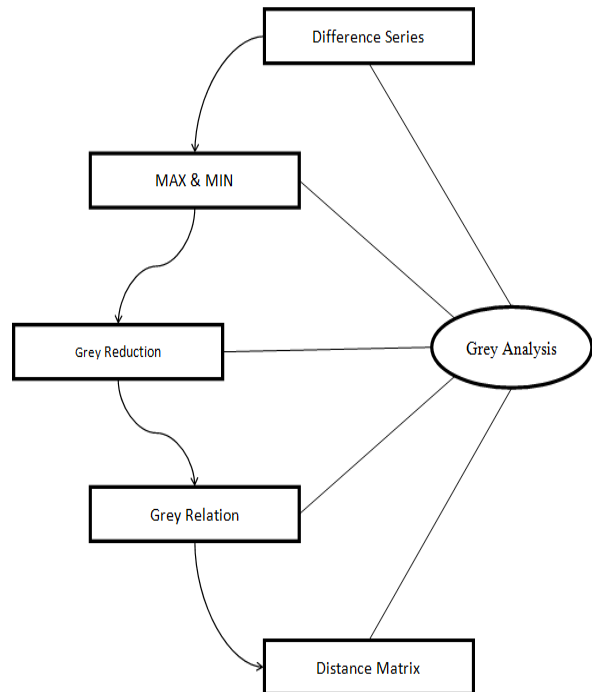


Figure 3.2: Process Diagram for Grey Clustering

This model is often applied for predicting decision making in industrial engineering and management science. Grey Relation System analyzes the impact of change between two events or components and is a simple decision process technique that is described in Figure 2.

ALGORITHM 2.0 GRA

Begin

$x = \{x_1, x_2, x_3 \dots x_n\}$ //Select Standard Vector.
 $\Delta(k) = |x_i(k) - x_j(k)|$ //Difference Matrix
 $\Delta_{max[i,j]} = \max |x_i(k) - x_j(k)|$ and $\Delta_{min[i,j]} = \min |x_i(k) - x_j(k)|$
 $Gr(x_i(k), x_j(k)) = (\Delta_{min} + \delta \Delta_{max}) \div (\Delta_{ij}(k) + \delta \Delta_{max})$ //
 $\varpi(i, j) = \frac{1}{k} \sum_{k=1}^k Gr(x_i(k), x_j(k))$ // Grey Relation
 $G(i, j) = \frac{(\varpi(i, j) + \max(\varpi(i, j)))}{2 \max(\varpi(i, j))}$

End

In respect to analysis we are applying Rough Set theory for reducing its superfluous attribute and for adjusting the affine objects we analyze the discretized data after RST approach. In our next step, we used Gray Relation method for clustering the similar objects. The main obstacles facing current Data Analysis techniques are that of dataset dimensionality. Usually, a redundancy-removing step is carried out beforehand to enable these techniques to be effective. Rough Set Theory (RST) has been used as such a dataset pre-processor with much success, however it is reliant upon a discretized dataset; but in some case the important information may be lost as a result of discretisation.

Step 2: If there any dispensable attribute in data set then Reduct otherwise make it as Indispensable elements.

4. EXPERIMENT

In this paper, we are using an ecological data set of plant diversity of North American Island, which consists lots of attributes that can affect the Richness of plant.

4.1 Rule Generation

Rule generation will generate the rules based on reduct and core of Table 2. It's produced the reduced set Rough-Set of relation that can transform the same inductive classification of Relation.

Table 4.1: Plant Diversity Data-Set [*]

	Island	tot.ri ch	ntv.ri ch	no.ri ch	Pct	Area	latitude	ele v	dist.m nland	dist.islan d	Soil Type
1	Appledore	182	79	103	57	40	42.99	18	10	10	6
2	Bear	64	43	21	33	3	41.25	13	0.3	0.3	1
3	Block	661	396	265	40	2707	41.18	64	20.6	20.6	59
4	Cuttyhunk	311	173	138	44	61	41.42	46	10.8	0.4	11
5	Fishers	920	516	404	44	1190	41.27	40	2.7	2.7	35
6	Gardiners	390	249	141	36	1350	41.08	37	6.7	6.7	37
7	Grand Ma.	633	374	259	41	13600	44.75	122	17.5	17.5	.
8	Gull Rock	34	15	19	56	4	44.96	10	13.2	1	.
9	Horse	107	75	32	30	4	41.24	10	1.9	0.3	1
10	Isle au Haut	641	370	271	42	1900	44.05	165	22.9	8.1	21
11	Kent Island	232	120	112	48	128	44.58	20	30.1	7	.
12	Machias S.	72	24	48	67	10	44.5	6	17.7	17.7	.
13	Martha's V	979	605	374	38	13600	41.39	95	13.4	13.4	47
14	Matinicus	62	21	41	66	8	43.79	15	30.6	4.7	1
15	Mount	1060	620	440	42	26668	44.33	466	0.3	0.3	74
16	Muskeget	156	88	68	44	140	41.33	10	35.7	7.5	4
17	Nantucket	1166	625	541	46	10900	41.27	33	42.5	21	27
18	Naushon	564	362	202	36	2300	41.47	53	8.6	8.6	18
19	Penikese	347	181	166	48	34	41.45	21	8.5	1.6	6
20	Tuckernuck	353	224	129	37	350	41.3	15	34	3	16
21	Whaleboat	163	99	64	39	47	43.76	23	1.3	1.3	4
22	Wooden B.	155	69	86	55	46	43.86	19	27.4	4.3	2

The set P of attributes is the reduct (or covering) of another set Q of attributes if P is minimal and the indiscernibility relations, defined by P and Q are same.

$$core = \cap reduct$$

In applying Reduct method we eliminates the superfluous information from Table 1 and regenerate another table with having those attribute which are more better associate with other values.

4.2 Grey Relation Analysis of Reduction Table

The hierarchical grey relation clustering analysis calculation has been process in following steps:

Step 4: Calculate the difference of values:

Where $\Delta_{ij}(k)$ difference function and $x_i(k)$ represents the i and j row respectively.

$$\Delta_{ij}(k) = |x_i(k) - x_j(k)|$$

Where $i, j \in \{1,2,3,4 \dots n\}$ and $k = \{1,2\}$

Table 4.2: Reduct Table

Total Rich	Area	Elevation	Non-native_richness
(155,329]	(41.5,44]		Inf]
[-Inf,155]	[-Inf,41.3]	[-Inf,7.15]	[-Inf,38.2]
(639, Inf]	[-Inf,41.3]	(13.3,26.3]	(38.2,43]
(155,329]	(41.3,41.5]	(7.15,13.3]	(43,48]
(639, Inf]	[-Inf,41.3]	[-Inf,7.15]	(43,48]
(329,639]	[-Inf,41.3]	[-Inf,7.15]	[-Inf,38.2]
(329,639]	(44, Inf]	(13.3,26.3]	(38.2,43]
(155,329]	(44, Inf]	(7.15,13.3]	(48, Inf]
[-Inf,155]	[-Inf,41.3]	[-Inf,7.15]	[-Inf,38.2]
(639, Inf]	(44, Inf]	(13.3,26.3]	(38.2,43]
(155,329]	(44, Inf]	(26.3, Inf]	(48, Inf]
[-Inf,155]	(44, Inf]	(13.3,26.3]	[-Inf,38.2]
(639, Inf]	(41.3,41.5]	(13.3,26.3]	(38.2,43]
[-Inf,155]	(41.5,44]	(26.3, Inf]	(43,48]
(639, Inf]	(44, Inf]	[-Inf,7.15]	(48, Inf]
(155,329]	(41.3,41.5]	(26.3, Inf]	[-Inf,38.2]
(639, Inf]	[-Inf,41.3]	(26.3, Inf]	(48, Inf]
(329,639]	(41.5,44]	(7.15,13.3]	(38.2,43]
(329,639]	(41.3,41.5]	(7.15,13.3]	(43,48]
(329,639]	(41.3,41.5]	(26.3, Inf]	(43,48]
(155,329]	(41.5,44]	[-Inf,7.15]	[-Inf,38.2]
[-Inf,155]	(41.5,44]	(26.3, Inf]	(43,48]

Table 4.3: Reduct Plant Diversity Data

ID	American Island	Plant Richness (Diversity)	Ground Elevation
1	Appledore Island	182	18
2	Bear Island	64	13
3	Block Island	661	64
4	Cuttyhunk Island	311	46
5	Fishers Island	920	40
6	Gardiners Island	390	37
7	Grand Manan Island	633	122

Table 4.4: Difference Table

ID	American Island	Plant Richness (Diversity)	Ground Elevation
1	Appledore Island	0	0
2	Bear Island	118	5
3	Block Island	479	46
4	Cuttyhunk Island	129	28
5	Fishers Island	738	22
6	Gardiners Island	208	19
7	Grand Manan Island	451	104

Calculate the Maximum and Minimum values of the difference series.

$$\Delta_{max [i,j]} = \max |x_i(k) - x_j(k)| \text{ and } \Delta_{min [i,j]} = \min |x_i(k) - x_j(k)|$$

Calculate grey relation Coefficient

$$Gr(x_i(k), x_j(k)) = (\Delta_{min} + \delta \Delta_{max}) \div \Delta_{ij}(k) + \delta \Delta_{max}$$

Where $\delta = 0.1$ an adjustable variable, $i, j \in \{1,2,3,4 \dots n\}$ and $k = \{1,2\}$

Table 4.5: Grey Relation Table

ID	American Island	Plant Richness (Diversity)	Ground Elevation
1	Appledore Island	1	1
2	Bear Island	0.384	0.936
3	Block Island	0.133	0.616
4	Cuttyhunk Island	0.363	0.724
5	Fishers Island	0.090	0.770
6	Gardiners Island	0.261	0.795
7	Grand Manan Island	0.140	0.415

Step 6:

Where $i, j \in \{1,2,3,4 \dots n\}$ and $k = \{1,2\}$

Calculate grey relation grade:

$$\rho_{i,j} = 1 \setminus k \sum_{k=1}^k Gr(x_i(k), x_j(k))$$

Table 4.6: Grey grade relation Table

ID	Appledore Island	Bear Island	Block Island	Cuttyhunk Island	Fishers Island	Gardiners Island	Grand Manan Island
Appledore	1.000	0.684	0.337	0.502	0.449	0.469	0.238
Bear Island	0.660	1.000	0.314	0.422	0.425	0.413	0.216
Block	0.374	0.375	1.000	0.459	0.514	0.412	0.582
Cuttyhunk	0.543	0.489	0.456	1.000	0.528	0.627	0.289
Fishers	0.430	0.425	0.456	0.500	1.000	0.518	0.287
Gardiners	0.528	0.494	0.434	0.653	0.552	1.000	0.293
Grand Manan	0.277	0.284	0.593	0.301	0.739	0.563	1.000

Step 7: Develop matrix G

$$G_{i,j} = (\rho_{i,j} + \rho_{j,i})/2$$

Table 4.7: Developed Grey Matrix G

ID	Appledore Island	Bear Island	Block Island	Cuttyhunk Island	Fishers Island	Gardiners Island	Grand Manan Island
Appledore	1.000						
Bear	0.672	1.000					
Block	0.355	0.344	1.000				
Cuttyhunk	0.522	0.455	0.457	1.000			
Fishers	0.439	0.425	0.485	0.514	1.000		
Gardiners	0.498	0.453	0.423	0.640	0.535	1.000	
Grand Manan	0.257	0.250	0.587	0.295	0.513	0.428	1.000

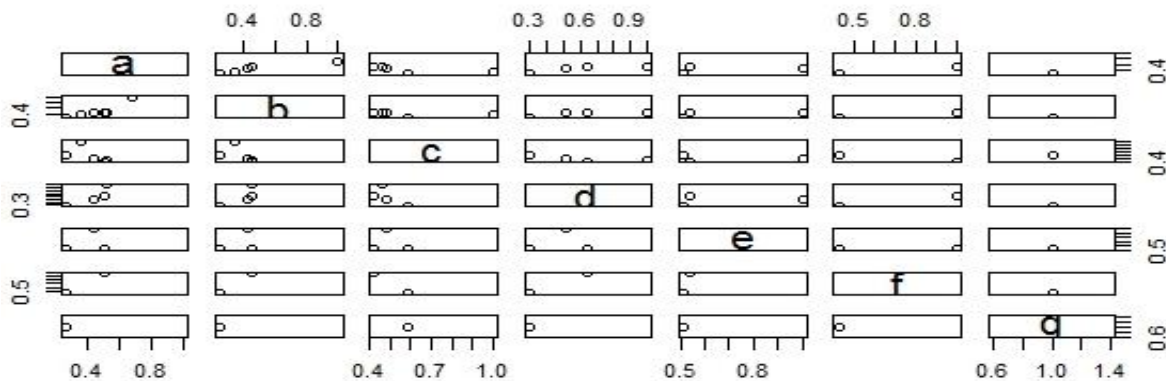


Fig 4.1: Relational Degree Plot

Step 8: Create Cluster by using comparison of two nearest point.

$$\max_{i,j}(G_{i,j})$$

Table 4.8: Clustering Table

Cluster-1 (0.500-0.600)	Gardiners Island, Bear Island, Appledore Island, Cuttyhunk Island, Fishers, Grand Manan Island.
Cluster-2(0.300-0.400)	Bear Island
Cluster-3(0.200-0.300)	Bear Island

Cluster Dendrogram

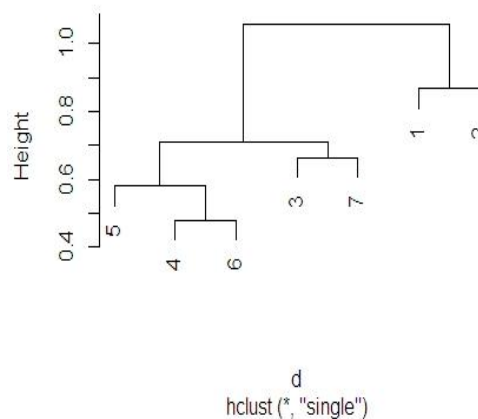


Fig. 4.2: Dendrogram Representation

5. CONCLUSION

This paper presents a new approach for extracting knowledge from large set of information, including high dimensional object analysis using Rough Set attribute reduction technique and using Grey Relational Clustering. I have used this data for ecological data set but have not explored its other applications. I have yet to compare this approach with other existing approaches. I expect this clustering approach to have benefits in data mining, agriculture, financial data analysis, biology, and several other fields.

7. Pawlak, Z., Rough sets. International Journal of Computer & Information Sciences, 1982. **11**(5): p. 341-356.
8. Maji, P., A.R. Roy, and R. Biswas, An application of soft sets in a decision making problem. Computers & Mathematics with Applications, 2002. **44**(8): p. 1077-1083.
9. Julong, D., Introduction to grey system theory. The Journal of grey system, 1989. **1**(1): p. 1-24.

ACKNOWLEDGEMENT

Our thanks to the experts who have contributed towards development of the template.

REFERENCES

1. Liu, X. and M. Li, Integrated constraint based clustering algorithm for high dimensional data. Neurocomputing, 2014. **142**: p. 478-485.
2. Aggarwal, C.C., et al. Fast algorithms for projected clustering. in ACM SIGMOD Record. 1999. ACM.
3. Burges, C.J., Dimension reduction: A guided tour. 2010: Now Publishers Inc.
4. Steinbach, M., L. Ertöz, and V. Kumar, The challenges of clustering high dimensional data, in New Directions in Statistical Physics. 2004, Springer. p. 273-309.
5. Verbeek, J., Mixture models for clustering and dimension reduction. 2004, Universiteit van Amsterdam.
6. Parsons, L., E. Haque, and H. Liu, Subspace clustering for high dimensional data: a review. ACM SIGKDD Explorations Newsletter, 2004. **6**(1): p. 90-105.