



Performance Analysis of Quickreduct, Quick Relative Reduct Algorithm and a New Proposed Algorithm

By Ashima Gawar

GGSIPU, India

Abstract- Feature Selection is a process of selecting a subset of relevant features from a huge dataset that satisfy method dependent criteria and thus minimize the cardinality and ensure that the accuracy and precision is not affected ,hence approximating the original class distribution of data from a given set of selected features. Feature selection and feature extraction are the two problems that we face when we want to select the best and important attributes from a given dataset Feature selection is a step in data mining that is done prior to other steps and is found to be very useful and effective in removing unimportant attributes so that the storage efficiency and accuracy of the dataset can be increased. From a huge pool of data available we want to extract useful and relevant information. The problem is not the unavailability of data , it is the quality of data that we lack in.. We have Rough Sets Theory which is very useful in extracting relevant attributes and help to increase the importance of the information system we have. Rough set theory works on the principle of classifying similar objects into classes with respect to some features and those features may collectively and shortly be termed as reducts.

Keywords: *data mining, rough set, quickreduct, quick relative reduct, feature selection, feature extraction.*

GJCST-C Classification: *F.2.0, I.1.2*



Strictly as per the compliance and regulations of:



Performance Analysis of Quickreduct, Quick Relative Reduct Algorithm and a New Proposed Algorithm

Ashima Gawar

Abstract- Feature Selection is a process of selecting a subset of relevant features from a huge dataset that satisfy method dependent criteria and thus minimize the cardinality and ensure that the accuracy and precision is not affected ,hence approximating the original class distribution of data from a given set of selected features. Feature selection and feature extraction are the two problems that we face when we want to select the best and important attributes from a given dataset Feature selection is a step in data mining that is done prior to other steps and is found to be very useful and effective in removing unimportant attributes so that the storage efficiency and accuracy of the dataset can be increased. From a huge pool of data available we want to extract useful and relevant information. The problem is not the unavailability of data, it is the quality of data that we lack in. We have Rough Sets Theory which is very useful in extracting relevant attributes and help to increase the importance of the information system we have. Rough set theory works on the principle of classifying similar objects into classes with respect to some features and those features may collectively and shortly be termed as reducts.

In this paper, we have discussed Quickreduct and Quick Relative Reduct algorithm and also proposed a new algorithm. A comparative study between these two algorithms is also done. The experimental results show that Quick Relative Reduct algorithm is better than Quickreduct algorithm. The analysis is carried out on synthetic datasets.

Keywords: data mining, rough set, quickreduct, quick relative reduct, feature selection, feature extraction.

I. INTRODUCTION

a) Feature Selection and Feature Extraction

It to the process of finding out and select minimum subsets of attribute from a large set of original attributes and finally select the minimal one. The aim behind the process is to reduce dimensions across the datasets, remove the attributes which have no significance and identify the most important and useful attributes. (Zhang et al., 2003) It will help in improving and increasing accuracy and lessen the time that the algorithm will take for its computation.

We have organized the remaining paper as follows : section 2 briefs about the data set used for the study. Section 3 describes the Quickreduct algorithm. Section 4 describes Quick Relative Reduct algorithm.

Author: MCA from Institute of Information Technology and Management, Janakpuri, New Delhi. e-mail: ashima_gavar@yahoo.in

Section 5 explains the analysis of the comparison made between the Quickreduct and the Quick Relative Reduct algorithm. Section 6 suggests some improvement in the QuickReduct algorithm and finally Section 7 states the conclusion of the paper.

b) Reducts

The minimal set of attributes that will identify the other attributes of the dataset thus improving its accuracy and efficiency are called reducts.(Jothi and Inbarani , October 2012) Mathematically , a reduct of an algebraic structure that is calculated by removing some of the operations and relations of the mathematical structure we are using. In a reduct we keep only those attributes that are similar in nature and consequently have the goal of set approximation . Usually we can find several such subsets and those which are minimal among those are called reducts.

Given an information table S, an attribute set $R \subseteq At$ is called a reduct, if R satisfies the two conditions:

1. $IND(R) = IND(At)$;
2. For any $a \in R$, $IND(R - \{a\}) \neq IND(At)$.

c) Rough Sets

Rough set theory provide a novel methodological approach for approximation of large sets and describing the knowledge. In rough set theory firstly we collect a sample object set and store the feature values in information tables. Rough sets help us to find reducts without deteriorating the original quality of the dataset.

Characterization of Rough sets cannot be done in terms of information about the elements of rough sets. With every rough set a pair of precise sets, known as the lower and the upper approximations of the rough set. The lower approximation contains all the objects which definitely belong to the set and the upper approximation contains all objects which may possibly belong to the set. The difference between the upper and the lower approximation constitutes the boundary region of the rough set. Approximations are the fundamental concepts of rough set theory.

Rough set theory can be described as a formal methodology that can be employed to reduce the dimensions of datasets and is used as an preprocessing step to data mining. The reduced

dimensionality improves the runtime performance of an algorithm. Rough Set theory (Suguna and Thanushkodi , 2010) is a mathematical approach that is based on the principle that if the degree of precision in a dataset is lowered then we can more easily visualize the data patterns. The main aim is to approximate the lower and upper bounds. Rough set based data analysis initially analyses the data table called decision table in which the columns are labeled by attributes and rows represent the objects. The entries of the table will contain the value of the attributes . Attributes of the decision table are divided into two disjoint groups which are called decision and condition attributes respectively. Any rough set is associated with a pair of precise sets which are called the lower and upper approximations of the rough set is associated (Yiyu and Yan 2009).

II. DATA PREPARATION

We have manually performed analysis on the test datasets . The first dataset contains information about AUTOMOBILE and the second one contains information contains data about COMPUTER.

III. QUICKREDUCT ALGORITHM (QR)

In Quickreduct algorithm we remove the attributes so that the set we get after reduction provides

the same prediction of the decision feature as the original set which is achieved by comparing equivalence relations generated by sets of attributes. The attribute selected for the first time is to be included in the reduct set in the Quickreduct algorithm (Velayutham and Thangavel, September 2011) is the degree of dependency of that attribute which is not equal to zero..The algorithm tries to find out a minimal reduct without generating all possible subsets . Initially we take an empty set and add in the empty set R those attributes that will result in the greatest increase in dependency value one by one until we get the maximum possible value for the dataset.

Algorithm

1. $R \leftarrow \{ \}$ // empty set
2. Do
3. $T \leftarrow R$ // take T
4. For all $x \in (C-R)$ // C is the core
5. If $Y_{R \cup \{x\}}(D) > Y_T(D)$
6. $T \leftarrow R \cup \{x\}$
7. $R \leftarrow T$
8. Until $Y_R(D) == Y_C(D)$
9. Return R

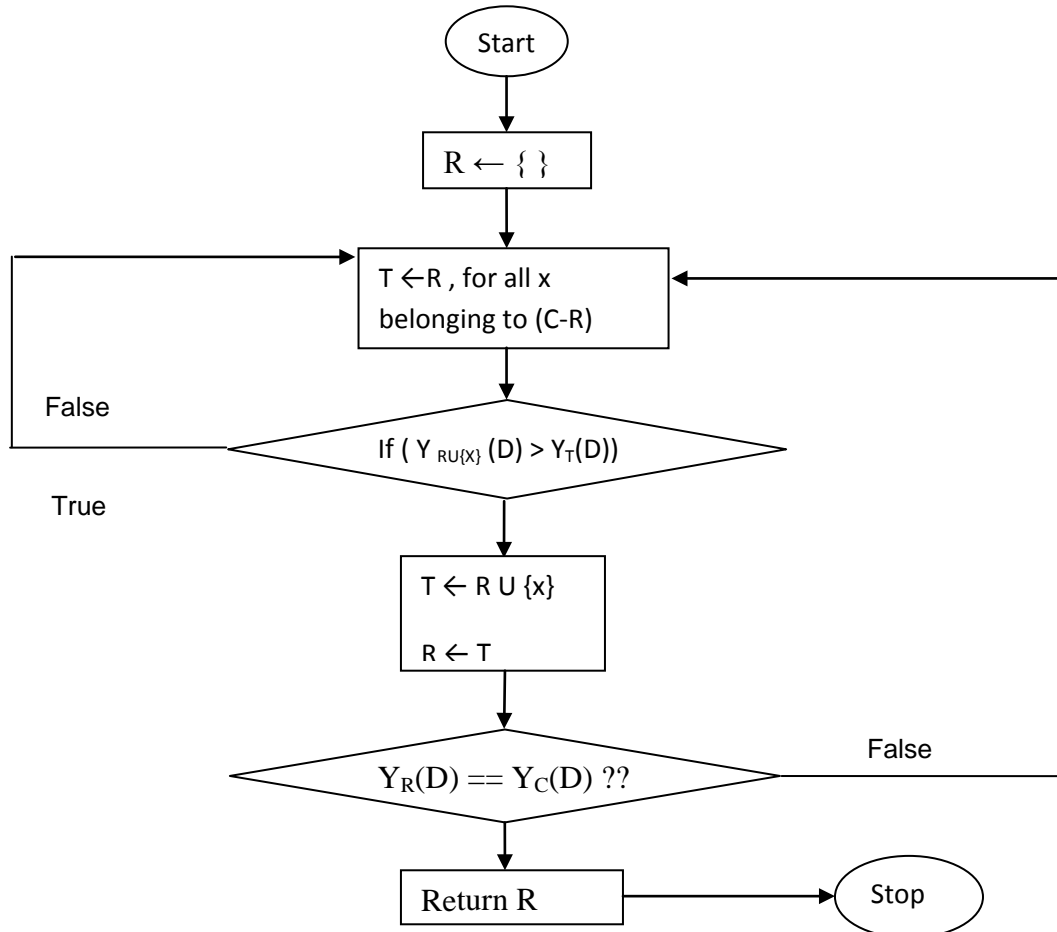


Figure 1: Stepwise execution of Quickreduct algorithm

IV. QUICK RELATIVE REDUCT ALGORITHM

In Quick Relative Reduct (Kalyani and karnan 2011)algorithm we find out the degree of relative dependency after removing the attributes from the set. If a attribute is removed and it causes the value of relative dependency to be one then that attribute is eliminated otherwise it is put in the core reduct. The process is repeated again and again till the value becomes one. The algorithm is explained below :

The algorithm is explained below :

- Algorithm
1. $R \leftarrow \{ \text{list of conditional attributes} \}$
 2. Select $x \leftarrow$ Conditional attribute from R
 3. If dependency = 1
 4. Then "Eliminate the attribute "
 5. Else
 6. $R \leftarrow \{ \text{list of reduct} \}$
 7. End

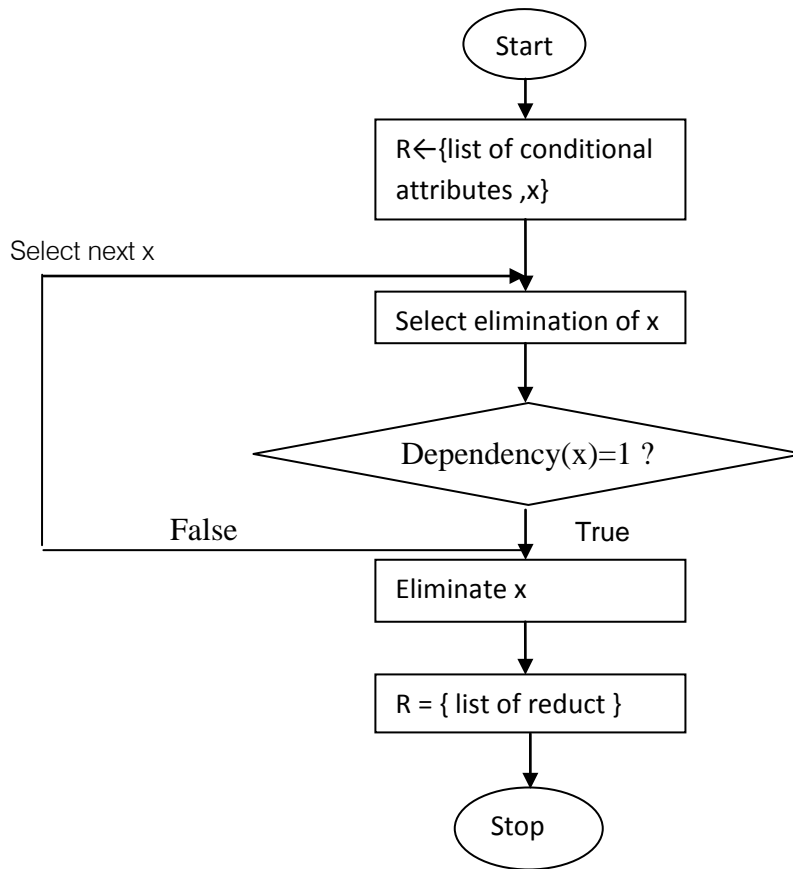


Figure 2 : Stepwise execution of Quick Relative Reduct algorithm

V. PERFORMANCE ANALYSIS

Table 1 : Results for Quick Reduct

Date set	Attributes	Instances	Selected Attributes	Reduct ?	Optimal ?
Automobile	4	8	3	Yes	Yes
Computer	6	20	4	Yes	No

Table 2 : Results for Quick Relative Reduct

Date set	Attributes	Instances	Selected Attributes	Reduct ?	Optimal ?
Automobile	4	8	3	Yes	Yes
Computer	6	20	3	Yes	Yes

Both Quick Reduct and Quick Relative reduct are reduct algorithms but the Quick Relative Reduct is a more efficient algorithm as it calculates reducts without

calculating discernibility functions which can be expensive. It includes a simple approach using relative dependency.

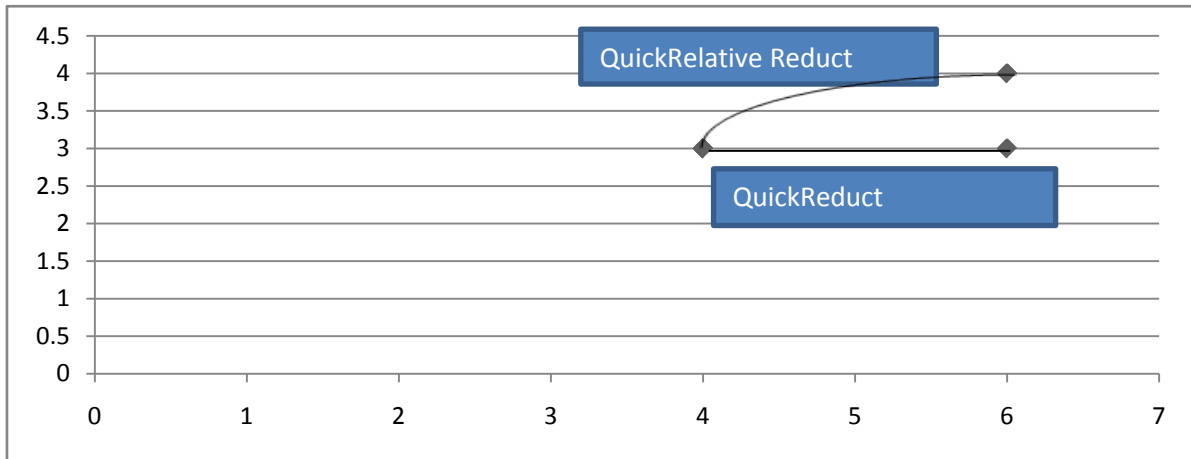


Figure 3 : Graph depicting performance analysis

VI. PROPOSED ALGORITHM

We propose a new algorithm to overcome the disadvantage of Quick Relative Reduct that in this algorithm we calculate relative dependency and the attribute is chosen with highest degree of dependency. When the highest relative dependency value is possessed by more than one attribute. For that purpose we can introduce a significance factor associated with every attribute and choose the attribute with greater significance. Significance factor (Jothi and Inbarani 2012) is defined as :

Assume $X \subseteq A$ is an attribute subset, $x \in A$ is an attribute, the importance of x for X is denoted by $Sig X(x) = 1 - \frac{|X \cup \{x\}|}{|X|}$ Where $|X| = |IND(X)|$. Suppose $U/IND(X) = U/X = \{X_1, X_2 \dots X_n\}$, then $|X| = |IND(X)| = \sum |X_i|$. $|X| - |X \cup \{x\}|$ represents the decrement of indiscernibility and also the increment of discernibility as attribute x is added to X . The number of selection methods is originally indiscernible in X , but it is discernible in $X \cup \{x\}$ and the increment of indiscernibility is expressed by :

$$\frac{(|X| - |X \cup \{x\}|)}{|X|} = 1 - \frac{|X \cup \{x\}|}{|X|} \quad (6.1)$$

Proposed Algorithm :

Input : x : conditional attributes

Step 1 : Take the R as the set of all conditional attributes.

Step 2 : Now select the conditional attribute.

Step 3 : Calculate the relative dependency of the attribute.

Step 4 : If relative dependency of the attribute is one then eliminate the attribute , Go to step 2.

Step 5 : If relative dependency is not equal to one then select the highest dependency attribute ,if two attributes

have same relative dependency then select the one with greater significance.

Step 6 : $R = \{ \text{list of reduct} \}$

Step 7 : Stop

Output : Reduct

VII. CONCLUSION

In this we discussed the comparison analysis of the Quickreduct and the Relative QuickReduct algorithm. The Relative QuickReduct algorithm finds reducts based on backward elimination of attributes and the QuickReduct algorithm finds reducts based on forward elimination. We also found out that Quick Relative Reduct was better than the QuickReduct algorithm. Also the Relative QuickReduct algorithm can be modified further to improve the efficiency by introducing the concept of significance factor. Further work can be carried out on the defined algorithm to explore its efficiency and accuracy. The analysis was performed manually but the research can be carried out further for further suggestions and improvements.

REFERENCES RÉFÉRENCES REFERENCIAS

1. Jothi. G and Inbarani H., Soft Set Based Feature Selection Approach for Lung Cancer Images , International Journal of Scientific and Engineering Research (October 2012) , Volume 3 Issue 10 , pp.1293-1299.
2. Z.Pawlak, Rough sets, International Journal of Computer Information and Science 11 (1982) pp. 341-356.
3. Velayutham C. and Thangavel K , Unsupervised Quick Reduct Algorithm Using Rough Set Theory , Journal Of Electronic Science And Technology (September 2011), VOL. 9, NO. 3 , pp. 193-201.

4. Zhang J. ,Wang J. , Huacan H., and Jianguang S. , A New Heuristic Reduct Algorithm Based on Rough Sets Theory , 4th International Conference, WAIM , Chengdu, China, August 17-19, 2003 , Volume 2762, 2003, pp 247-253.
5. Suguna N. and Dr. Thanushkodi K. , A Novel Rough Set Reduct Algorithm for Medical Domain Based on Bee Colony Optimization , Journal Of Computing, Volume 2, Issue 6, June 2010 , pp. 49-54.
6. Chandrasekhar T. , Thangavel K. and Sathishkumar E.N . , Verdict Accuracy of Quick Reduct Algorithm using Clustering and Classification Techniques for Gene Expression Data , IJCSI International Journal of Computer Science Issues ,January 2012 , Vol. 9, Issue 1, No1 , pp. 357-363.
7. Slowinski, R., Vanderpooten, D., A generalized definition of rough approximations based on similarity (2000) , IEEE Transactions on Knowledge and Data Engineering, pp. 331-336.
8. Yiyu Y and Yan Z , Discernibility Matrix Simplification for Constructing Attribute Reducts Discernibility matrix simplification for constructing attribute reducts, Information Sciences, 2009 Vol. 179, No. 5, pp. 867-882.
9. Yao, Y.Y., Wong, S.K.M, Lin, T.Y., A review of rough set models, (1997) Rough Sets and Data Mining: Analysis for Imprecise Data, pp. 47-73.



This page is intentionally left blank