ELSEVIER

2013 AASRI Conference on Intelligent Systems and Control

# Clustering Software Components for Program Restructuring and Component Reuse Using Hybrid XOR Similarity Function

Chintakindi Srinivas,Vangipuram Radhakrishna, Dr.C.V.Guru Rao

*Associate Professor of CSE, Kakatiya Institute of Technology and Science, Warangal, INDIA*
*Department of Information Technology, VNR Vignana Jyothi Institute of Engineering and Technolog, Hyderabad,INDIA*
*Professor of Computer Science and Engineering, SR Engineering College,Warangal,INDIA*

**Abstract**

Component based software development has gained a lot of practical importance in the field of software engineering from several researchers and also from industry perspective. Finding components for efficient software reuse is one of the important problems aimed by researchers. Clustering reduces the search space of components by grouping similar entities together thus ensuring reduced time complexity as it reduces the search time for component retrieval. In this research, we instigate a new a generalized approach for clustering a given set of documents or text files or components by defining a new similarity function called hybrid XOR function for the purpose of finding degree of similarity between two document sets or software components. We construct a matrix called similarity matrix of the order n-1 by n for a given set of n documents or components or patterns by applying hybrid XOR function. We define and design the algorithm for component or document clustering which has the input as similarity matrix and output being set of clusters formed dynamically as compared to other clustering algorithms that predefine the count of clusters. The output is a set of highly cohesive pattern groups or documents. The approach can be justified as it carries out very simple computational logic and efficient in terms of processing with reduced search space and can be also be used in general for document clustering or software component clustering.

*Keywords :* Clustering, hybrid XOR;frequent itemsets; mining;classification,components

## 1. Introduction

Clustering is one of the topics which have achieved a lot of practical importance from the researchers and also from the perspective of the software industry. The significance for clustering approach comes from the

---

\* Corresponding author. Tel.: +9700684242
*E-mail address:* radhakrishna_v@vnrvjiet.in

need of decision making such as classification, prediction , component search and retrieval and is thus widely used in many practical domains such as text classification , bioinformatics, medicine , image processing. We can define Clustering as the process of grouping similar set of patterns together [11]. The input to clustering algorithm may be any set of entities or patterns or text files or software components. The output of clustering algorithm will be a partition of cohesive groups. The abstract representation of clustering process is shown in figure 1 below.
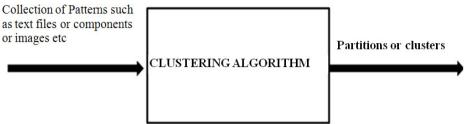


Figure 1:  Abstract view of clustering process

The representations or descriptions of clusters so formed shall be used in decision making such as which software component or pattern need to be selected. One important feature of clustering is that all the patterns within a cluster share common or same properties in some sense and patterns in different clusters are dissimilar in corresponding sense. In view of software engineering, all the components within the same cluster have high cohesion and low coupling.

Software component clusters can be treated as highly cohesive groups with low coupling which is the desired feature. One disadvantage of existing data clustering methods is that they do not adequately address the problem of processing large datasets with a limited amount of resources. Using these limitations as our motivation, so if we can try to reduce the dataset for training process it can help in reducing the cost of training which in turn improves efficiency of clustering. If done so, clustering takes less amount of space and hence forms a compact storage of patterns. Clustering is not any one specific algorithm that we can stick firm to, but it must be viewed as the general task to be solved.

Clustering algorithms may unsupervised or supervised [11]. In unsupervised clustering the partitions are viewed as the unlabelled patterns or components. Supervised clustering algorithms label the patterns which can be used to classify the components for decision making.  Hence the partitions obtained by clustering process may be labeled or unlabeled.

A new method called Maximum Capturing is proposed for document clustering [3].Maximum Capturing includes two procedures:  1. constructing document clusters and 2. Assigning cluster topics. The search complexity can be reduced by using the algorithm [10] where ever necessary as part of component retrieval.

## 2. Taxonomy

The problem of finding frequent itemsets gets birth from [8] which uses frequent itemsets to find association rules of items in large transactional databases. In  [1] clustering a given set of text documents from neighbour set is proposed. In [2] the authors propose a method for discovering maximum length frequent item sets. In [6], the classification of text files or documents is done by considering Gaussian membership function and making use of it to obtain clusters by finding word patterns. Each cluster is identified by its word pattern calculated using fuzzy based Gaussian membership function once clusters are formed.

In this paper the idea is to first obtain frequent item sets for each document using existing association rule mining algorithms either by horizontal or vertical approach. Once we find frequent itemsets in each document

then we form a Boolean matrix with rows indicating documents and columns indicating unique frequent item sets from each document. This is followed by the computation of a binary feature vector for each document pair, represented as a 2D array or 2D matrix by redefining the XOR function as hybrid XOR logic with slight modification in the function introducing high impedance variable as Z. The idea of maximum capturing is taken as the base framework for clustering.

## 3. Proposed Work

To design a clustering algorithm we must first design the similarity function which is the heart of any clustering algorithm. We define a generalized similarity function called Hybrid XOR function which may be used to compute similarity feature between any pair of entities which may be software components or software patterns or documents. The documents may be text files to be classified or software product documents of various phases in software life cycle. We define the similarity function S as a function of any two entities A and B which is a tri state function as shown below in the truth table 1.

Table 1. Truth Table of hybrid XOR Similarity Function

| A | B | S(A,B) |
|---|---|--------|
| 0 | 0 | Z |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

**Hypothesis-1:** If a frequent item set exists in the document, then the cell value of the matrix corresponding to

D $[d_i, w_k]$ is made 1 else the corresponding cell value of the matrix is made as zero.

The algorithm for document clustering has its input as documents with frequent item sets and output as set of clusters formed dynamically. The approach followed is a tabular approach. Similarly the algorithm for component clustering has its input as software components with properties predefined and the output is a set of highly cohesive components with low coupling feature.

### 3.1 Algorithm for Clustering

// may be used for software component clustering or document clustering or pattern clustering in general
Document_Clustering (Document set, frequent item sets)
Begin of Algorithm
Step1:For each document D do
 Begin
 Step1.   Remove stop words and stemming words from each document.
 Step2.   Find unique words in each document and count of the same.
 Step3.   Find frequent itemsets of each document
 End for
 Step 2: Form a word set W consisting of each word in frequent item sets of each document.
 Step 3: Form Dependency Boolean Matrix with each row and column corresponding to each Document and each word respectivelyFor each document in document set do
 Begin
 For each word in word set do

Begin
If (word $W_k$ in Word set W is in document $D_i$)
Begin
Set $D[D_i, W_k] = 1$
Else
Set $D[D_i, W_k] = 0$
End if
End for
End for

Step 4 : Find the Feature vector similarity matrix by evaluating similarity value for each document pair applying Hybrid XOR Function defined in table 1 to obtain the matrix with feature vectors for each document pair.

Step 5: Replace the corresponding cells of matrix by count of number of zeroes in tri state feature vector.

Step 6: At each step, find the cell with maximum value and document pairs containing this value in the matrix.Group such document pairs to form clusters. Also if document pair (I, J) is in one cluster and document pair (J, K) is in another cluster, form a new cluster containing (I, J, K) as its elements.

Step 7: Repeat Step6 until no documents exist or we reach the stage of first minimum value leaving zero entry.

Step 8: Output the set of clusters obtained.

Step 9: label the clusters by considering candidate entries.
*End of algorithm*

### 3.2 Case Study showing process of document clustering

Consider the document sets with the frequent item sets obtained after mining using any of the existing association rule mining algorithms as shown below. Here we use can also use association rule mining algorithm with multiple support and confidence thresholds. We considered a set of random of 20 documents as the training set.

Table 2. Documents and Corresponding Frequent item sets

| DOCUMENTS | FREQUENT ITEMSETS |
|---|---|
| DOCUMENT 1 | { ENCRYPT,  NEURAL NETWORKS,  CLUSTER} |
| DOCUMENT 2 | {SVM , MINING,  CLUSTER} |
| DOCUMENT 3 | { SVM, NEURAL NETWORKS, MINING,  CLUSTER} |
| DOCUMENT 4 | {ENCRYPT,  NEURAL NETWORKS, MINING,  CLUSTER} |
| DOCUMENT 5 | {SVM , CLUSTER} |
| DOCUMENT 6 | {ENCRYPT,   NEURAL NETWORKS, MINING} |
| DOCUMENT 7 | {ENCRYPT, SVM, NEURAL NETWORKS} |
| DOCUMENT 8 | {SVM, NEURAL NETWORKS} |
| DOCUMENT 9 | {NEURAL NETWORKS, MINING,  CLUSTER} |

We now construct a Boolean matrix with rows indicating each document and column corresponding to each unique frequent item from set of frequent item sets of all documents sets respectively.

Table 3. Boolean matrix Representation of Table.2

|  | ENCRYPT | NEURAL NETWORKS | CLUSTER | SVM | MINING |
|---|---|---|---|---|---|
| D1 | 1 | 1 | 1 | 0 | 0 |
| D2 | 0 | 0 | 1 | 1 | 1 |
| D3 | 0 | 1 | 1 | 1 | 1 |
| D4 | 1 | 1 | 1 | 0 | 1 |
| D5 | 0 | 0 | 1 | 1 | 0 |
| D6 | 1 | 1 | 0 | 0 | 1 |
| D7 | 1 | 1 | 0 | 1 | 0 |
| D8 | 0 | 1 | 0 | 1 | 0 |
| D9 | 0 | 1 | 1 | 0 | 1 |

We form a matrix D [n-1, n] for n documents and consider only the upper triangular region. The cells of the matrix are filled by applying the similarity function S for which each document pair forms the input as shown below

Table 4. Feature Vector Representation of document set

|  | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 |
|---|---|---|---|---|---|---|---|---|---|
| D1 | x | {1,1,0,1,1} = 1 | {1,0,0,1,1}=2 | {0,0,0,Z,1}=3 | {1,1,0,1,Z=1 | {0,0,1,Z,1}=2 | {0,0,1,1,Z}=2 | {1,0,1,1,Z}=1 | {1,0,0,Z,1=2 |
| D2 | x | x | {Z,1,0,0,0}=3 | {0,0,0,0,1}=4 | {Z,Z,0,0,=2 | {1,1,1,1,0}=1 | {1,1,1,0,1}=1 | {Z,1,1,0,1}=1 | {Z,1,0,1,0}=2 |
| D3 | x | x | x | {1,0,0,1,0}=3 | {Z,1,0,0,1}=2 | {1,0,1,1,0}=3 | {1,0,1,0,1}=2 | {Z,0,1,0,1}=2 | {Z,0,0,1,0}=3 |
| D4 | x | x | x | x | {1,1,0,1,1}=1 | {0,0,1,0,0}=4 | {0,0,1,1,1 }=2 | { 1,0,1,1,1}=1 | { 1,0,0,Z,0}=3 |
| D5 | x | x | x | x | x | {1,1,1,1,1}=0 | {1,1,1,0,Z}=1 | {Z,1,1,0,Z}=1 | {Z,1,0,1,1}=1 |
| D6 | x | x | x | x | x | x | {0,0,Z,1,1}=2 | {1,0,Z,1,1}=1 | {1,0,1,Z,0}=2 |
| D7 | x | x | x | x | x | x | x | {1,0,Z,0,Z}=2 | {1,0,1,1,1}=1 |
| D8 | x | x | x | x | x | x | x | x | {Z,0,1,1,1}=1 |

Once we obtain the above table with feature vectors for each document pair then we replace the corresponding cells of matrix by count of number of zeroes in tri state feature vector. We call it tri-state because it can have 0 or 1 or z as the value. This is shown in the table below.

Table 5. Similarity Matrix with Feature Vector Replaced by Count of 0s.

|  | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 |
|---|---|---|---|---|---|---|---|---|---|
| D1 | x | 1 | 2 | 3 | 1 | 2 | 2 | 1 | 2 |
| D2 | x | x | 3 | 4 | 2 | 1 | 1 | 1 | 2 |
| D3 | x | x | x | 3 | 2 | 3 | 2 | 2 | 3 |
| D4 | x | x | x | x | 1 | 4 | 2 | 1 | 3 |
| D5 | x | x | x | x | x | 0 | 1 | 1 | 1 |
| D6 | x | x | x | x | x | x | 2 | 1 | 2 |
| D7 | x | x | x | x | x | x | x | 2 | 1 |
| D8 | x | x | x | x | x | x | x | x | 1 |

Now consider only the element of the matrix with the highest value as shown in the table below. The step by step procedure is shown below in the form of tables which is self explanatory.

Step1: find the first maximum value from the matrix and target only those cells having this value to form initial cluster.

Table 6.  Content of Similarity Matrix showing step1

| | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | |
|---|---|---|---|---|---|---|---|---|---|---|
| D1 | x | 1 | 2 | 3 | 1 | 2 | 2 | 1 | 2 | **Find max value from the above table which is 4 here and target those cells as they form the best candidate solutions and replace those cell by x** |
| D2 | x | x | 3 | x | 2 | 1 | 1 | 1 | 2 | |
| D3 | x | x | x | 3 | 2 | 3 | 2 | 2 | 3 | |
| D4 | x | x | x | x | 1 | x | 2 | 1 | 3 | |
| D5 | x | x | x | x | x | 0 | 1 | 1 | 1 | **Stage1: (2, 4) and (4, 6) have val as 4. So form cluster as (2, 4, 6).** |
| D6 | x | x | x | x | x | x | 2 | 1 | 2 | |
| D7 | x | x | x | x | x | x | x | 2 | 1 | |
| D8 | x | x | x | x | x | x | x | x | 1 | |

Step 2: Find the next max value from the above table which is 3 here and target those cells as they form the best candidate solutions. Now cluster {2, 4, 6} is dynamically changed to {1, 2, 3, 4, 6, 9} and is no more a separate cluster as shown in table.2. cell values with superscript * not considered.

Table 7.  Content of Similarity Matrix showing step2

| | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | |
|---|---|---|---|---|---|---|---|---|---|---|
| D1 | x | 1 | 2 | 3 | 1 | 2 | 2 | 1 | 2 | **Find the next max value from the above table which is 3 here and target those cells as they form the best candidate solutions.** |
| D2 | x | x | 3* | x | 2 | 1 | 1 | 1 | 2 | |
| D3 | x | x | x | 3 | 2 | 3 | 2 | 2 | 3 | |
| D4 | x | x | x | x | 1 | x | 2 | 1 | 3* | **Stage2: consider only un-clustered document set {1, 3, 5, 7, 8, 9} ad search for value 3 in corresponding columns. (1,4)-(3,4)-(3,6)-(3,9) : So form cluster {2,4,6,1,3,9} as new Cluster. Set the values as zero or x.** |
| D5 | x | x | x | x | x | 0 | 1 | 1 | 1 | |
| D6 | x | x | x | x | x | x | 2 | 1 | 2 | |
| D7 | x | x | x | x | x | x | x | 2 | 1 | |
| D8 | x | x | x | x | x | x | x | x | 1 | **Cluster 1: {1, 2, 3, 4, 6, 9}** |

Step 3: Find the next max value from the above table which is 2 here and target those cells as they form the best candidate solutions.

Table 8.  Content of Similarity Matrix showing step3

| | D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | |
|---|---|---|---|---|---|---|---|---|---|---|
| D1 | x | 1 | 2 | x | 1 | 2 | 2 | 1 | 2 | **Find the next max value from the above table which is 2 here and target those cells as they form the best candidate solutions.** |
| D2 | x | x | x | x | 2 | 1 | 1 | 1 | 2 | |
| D3 | x | x | x | x | 2 | x | 2 | 2 | x | |
| D4 | x | x | x | x | 1 | x | 2 | 1 | x | **Stage3: consider only un-clustered document set {5, 7, 8} ad search for value 2 in corresponding columns. Here (7, 8) has 2. So form cluster {7, 8} as new Cluster. Set the values as zero or x.** |
| D5 | x | x | x | x | x | 0 | 1 | 1 | 1 | |
| D6 | x | x | x | x | x | x | 2 | 1 | 2 | |
| D7 | x | x | x | x | x | x | x | 2 | 1 | |
| D8 | x | x | x | x | x | x | x | x | 1 | **Cluster 2: {7, 8}** |

Step 4: Find the next max value from the above table which is 1 here and target those cells as they form the best candidate solutions.

Table 9.  Content of Similarity Matrix showing step4

| | d1 | d2 | d3 | d4 | d5 | d6 | d7 | d8 | d9 |
|---|---|---|---|---|---|---|---|---|---|
| **D1** | x | 1 | 2 | x | 1 | 2 | 2 | 1 | 2 |
| **D2** | x | x | x | x | 2 | 1 | 1 | 1 | 2 |
| **D3** | x | x | x | x | 2 | x | 2 | 2 | x |
| **D4** | x | x | x | x | 1 | x | 2 | 1 | x |
| **D5** | x | x | x | x | x | 0 | 1 | 1 | 1 |
| **D6** | x | x | x | x | x | x | 2 | 1 | 2 |
| **D7** | x | x | x | x | x | x | x | x | 1 |
| **D8** | x | x | x | x | x | x | x | x | 1 |

**Stage4: consider only un-clustered document set {5} and search for value 1 in corresponding columns. Here (5, 7), (5, 8), (5, 9) are all 1s. But this is next minimum value after zero if we consider initial table values before clustering. Hence 5 can't be similar to any of those documents and we must place it as a separate cluster {5}.**

The Set of clusters finally formed are as shown below in the following figure.

Cluster-1: {1, 2, 3, 4, 6, 9}

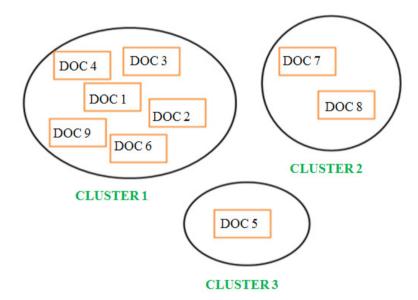Cluster-2 :{ 7, 8}

Cluster-3 :{ 5}



Figure. 2. Set of Clusters formed after applying the algorithm

## 4. Case Study for Software Component Clustering or Program Partitioning

Consider the following program fragment
```
Procedure Sum_and_Prod (n: integer; arr: int_array; var sum, prod: integer; var avg: float)
var i : integer;
begin
1. sum = 0;
2. prod = 1;
3. for i = 1 to n do begin
4. sum = sum + arr[i];
5. prod = prod * arr[i];
6. end;
```

7. avg = sum / n;
   end;

The table below shows the matrix with rows denoting line numbers or statements, columns denoting variable names.

Table 10.  Boolean Matrix for program module

| Line Numbers | Sum | Prod | N | Arr | Avg |
|---|---|---|---|---|---|
| S1 | 1 | 0 | 0 | 0 | 0 |
| S2 | 0 | 1 | 0 | 0 | 0 |
| S4 | 1 | 0 | 0 | 1 | 0 |
| S5 | 0 | 1 | 0 | 1 | 0 |
| S7 | 1 | 0 | 1 | 0 | 1 |

The table below shows the similarity matrix formed using algorithm. The below list of tables show the step by step process of forming clusters and final output of clusters and are self descriptive.

Table 11.  Trace of algorithm for stage 1

| | S1 | S2 | S4 | S5 | S7 | |
|---|---|---|---|---|---|---|
| S1 | X | 0 | 1 | 0 | 1 | Group S1 and S4 as (S1, S4) =1 and Mark as X. Mark Row elements of S4 by X.  This is done to reduce overlapping of variables. So (S1, S4) forms one cluster. |
| S2 | | X | 0 | 1 | 0 | |
| S4 | | | X | 1 | 1 | |
| S5 | | | | X | 0 | |
| S7 | | | | | X | |

Table 12.  Trace of algorithm for stage 2

| | S1 | S2 | S4 | S5 | S7 | |
|---|---|---|---|---|---|---|
| S1 | X | 0 | X | 0 | 1 | Again we have (S1, S7) =1. So group S1 and S7 into one cluster and mark the cell as X. As already S7 row is marked no need to do so. |
| S2 | | X | 0 | 1 | 0 | |
| S4 | | | X | X | X | Since (S1, S4) are similar and (S1, S7) are similar hence (S1, S4, S7) are all similar and placed into one cluster. |
| S5 | | | | X | 0 | |
| S7 | | | | | X | **i.e Cluster1 = {S1,S4,S7}** |

Table 13.  Trace of algorithm for stage 3

| | S1 | S2 | S4 | S5 | S7 | |
|---|---|---|---|---|---|---|
| S1 | X | 0 | X | 0 | X | Finally we have (S2, S5) =1. So group S2 and S5 into one cluster and mark the corresponding cell as X. As already S5 row is marked no need to do so. |
| S2 | | X | 0 | 1 | 0 | **i.e Cluster2 = {S2,S5}** |
| S4 | | | X | X | X | |
| S5 | | | | X | 0 | |
| S7 | | | | | X | |

The program module may finally be separated into two individual cluster modules which may run separately

**Cluster1**

1. sum = 0;
4. sum = sum + arr[i];
7. avg = sum / n;

**Cluster2**

2. prod = 0;
5. prod= prod* arr[i];

Figure. 3.   Program partitioning to run two modules in parallel

## 5. Conclusion

In this paper an attempt is made to study the problem of clustering software components for developing reuse library files and also the method of document clustering. An algorithm to cluster a set of given documents or text files or software components is designed which uses the new similarity function defined in this paper named hybrid XOR function defined for the purpose of finding degree of similarity among any two entities. The Proposed algorithm has the input as similarity matrix and the output being set of clusters formed dynamically as compared to other clustering algorithms that predefine the count of clusters and documents being fit to one of those clusters or classes finally. The approach can be extended to classify using classifiers and applying fuzzy logic in future. The concept of Support vector machines may be used for classification once clusters are formed if required. The search complexity can be reduced by using the algorithm [10] where ever necessary as part of component retrieval.

## References

[1]   Congnan Luo, , Yanjun Li, Soon M. Chung. Text document clustering based on neighbors , Data & Knowledge Engineering (68), 2009, 1271–1288.
[2]   Tianming Hu,Sam Yuan Sung, Hui Xiong, Qian Fu. Discovery of maximum length frequent itemsets, Information Sciences   (178), 2008,69–87.
[3]   Wen Zhanga,, Taketoshi Yoshida, Xijin Tang, Qing Wang. Text clustering using frequent itemsets, Knowledge-Based Systems 23 (2010) 379–388
[4]   Wen Zhanga,Taketoshi Yoshida, Xijin Tang. A comparative study of TF*IDF, LSI and multi-words for text classification. Expert Systems  with Applications 38 (2011) 2758–2765.
[5]   Vincent Labatut and Hocine Cherifi. Accuracy Measures for the Comparison of Classifiers, ICIT 2011 The 5th International Conference on  Information Technology.
[6]   Jung-Yi Jiang et.al A Fuzzy Self-Constructing Feature Clustering  Algorithm for TextClassification, IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 23, NO. 3, MARCH 2011.
[7]   Melita Hajdinjak, Andrej Bauer. Similarity Measures for Relational Databases, Informatica 33 (2009) 143–149.

[8]  R. Agrawal, T. Imielinski, A. Swami. Mining association rules between sets of items in very large databases, Proceedings of the ACM SIGMOD Conference on Management of data, 1993, pp. 207–216.

[9]  F. Beil, M. Ester, X.W. Xu, Frequent term-based text clustering, in: Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2002, pp. 436–442.

[10] Radhakrishna.V,C.Srinivas, C.V.Guru rao.  High Performance Pattern Search algorithm using three sliding windows , International Journal of Computer Engineering and Technology , Volume 3,issue 2, 2012 , pages 543-552. Impact factor 3.85.

[11] V.Susheela Devi, M. Narasimha Murthy. Text Book on Pattern Recognition. An Introduction. University Press.

[12] Salim kebir, Abdelhak-djamel seriai, Sylvain Chardigny. Comparing and Combining Genetic and Cluster Algorithms for Software Component Identification, in the Proceedings of the ACM Fifth International Conference on Computer Science and Software Engineering, Pages 1-8, 2012.

[13] Ronaldo.C.Veras, Silvio R.L.Meira , Adriano L.I. Oliveira , Bruno J.M.Melo. Comparitive study of clustering techniques for the organisation of software repositories, in the proceedings of 19th IEEE International Conference on Tools with Artificial Intelligence.

Dr.  C.V.GuruRao is currently the Head of the Department of CSE at S.R.Engineering College, Warangal, Andhra Pradesh, India. He has nearly 30 Years of teaching experience. He is a double post graduate, with specializations in Electronic Instrumentation and Information Science & Engineering. He is Doctorate holder in Computer Science & Engineering from Indian Institute of Technology, Kharagpur, India. He has more than 30 publications to his credit. He also served as the Chairman, Board of Studies for Computer Science & Engineering and Information Technology, Kakatiya University, Warangal. He is also serving as the Editorial Board member for International Journal of Computational Intelligence Research and Application journal. He is a life member of Indian Society for Technical Education, Instrumentation Society of India, and member of Institution of Engineers, Institution of Electronics & Telecommunications Engineers and Institution of Electrical & Electronics Engineers (USA).

C.Srinivas is a Masters Degree holder in Computer Science and Engineering from JNTU, Hyderabad. Presently he is working as an Associate Professor in CSE Department at Kakatiya Institute of Technology and Sciences, Warangal and is a research scholar at Kakatiya University under the guidance of Dr.C.V.Guru Rao. He has over 15years of teaching experience and presented papers at several national and international conferences and workshops. His areas of interest are Software Reuse, Cloud Computing.