# New Approach for Joint Multilabel Classification with Community-Aware Label Graph Learning Technique

R. V. Argiddi
Assistant Professor
WIT College of Engg.
Department of Computer Science
Solapur-413006, India

Disha Rajan Shah
P.G. Student
WIT College of Engg.
WIT Department of Computer Science, WIT
Solapur-413006, India

## ABSTRACT

Multi-label classification is a significant machine learning task in which one allocates a subset of candidate labels to an object. A new multi-label classification technique based on Conditional Bernoulli Mixtures. Exploiting label dependency for multi-label image classification cans considerably develop classification performance. Probabilistic Graphical Models are one of the primary methods for demonstrating such dependences. The structure of graphical models, however, is which ever resolute heuristically or learned from very inadequate information. Moreover, neither of these methodologies scales well to large or complex graphs. We recommend a principled way to learn the structure of a graphical model by in view of input features and labels, composed with loss functions. We formulate this problem into a max-margin framework primarily, and then convert it into a convex programming problem. In conclusion, we suggest a highly scalable technique that activates a set of cliques iteratively. Our methodology exhibits both strong theoretical properties and a substantial performance development over state-of-the-art approaches on both synthetic and real-world data sets. Our proposed system has numerous attractive properties: it captures label dependences; it decreases the multi-label problem to numerous standard binary and multi-class problems; it subsumes the classic independent binary prediction and power-set subset prediction approaches as special cases; and it exhibitions accuracy and/or computational complexity benefits over present approaches. We demonstrate two implementations of our technique by means of logistic regressions and gradient boosted trees, organized with a simple training procedure centered on Expectation Maximization. We promote derive an efficient prediction procedure centered on dynamic programming, thus avoiding the cost of scrutinizing an exponential number of probable label subsets. For the testing we will use and show the efficiency of the proposed method in contradiction of competitive substitutes on benchmark datasets with image as well as pdf.

## Keywords

Multi graph learning, Document classification approach, semi supervised learning.

## 1. INTRODUCTION

Text classification is to map the text to one or additional pre-defined kinds by means of a kind of classification algorithm which is accomplished permitting to text content. A standard classification corpus has been established and a unified evaluation method is accepted to organize English text – grounded on machine learning which has made a large growth now. Most real-world data are stored in relational databases. So to categorize objects in one relation, other relations provide crucial information. Traditional mechanism cannot convert relational data into a single table without expert knowledge or loosing crucial information. Multi-relational classification automatically classifies objects using multiple relations. Massive amounts of real world data are regularly collected into and planned in relational databases. Greatest of today's structured data is stored in relational databases. Thus, the task of learning from relational data has begun to receive Noteworthycourtesy in the literature. Unfortunately, furthermost methods simply utilize flat data representations. Hence, to apply these single-table data mining techniques, it forces to sustain a computational fine to first transforming the data into this flat form. Patterns of activity that, in isolation, are of limited significance for classification but, when combined/related, will expand the performance of system. Multi – relational classification goals at find outbeneficial patterns from corner to corner multiple inter-connected tables (relations) in a relational database. Traditional machine learning methodsundertake a random sample of homogeneous data from single relation but real world data sets are multi-relational and heterogeneous. Current solution does not scale well and cannot realistically be applied when considering database containing huge amount of data.

## 2. LITERATURE REVIEW

**1] : Xi Li at. Al. Joint Multilabel Classification With Community-Aware Label Graph Learning in IEEE 2016**.
propose a multi label classification framework based on a joint learning method called label graph learning (LGL) driven weighted Support Vector Machine (SVM). In belief, the joint learning method explicitly models the inter-label correlations by LGL, which is jointly optimized with multi label cataloging in a unified learning scheme. As aoutcome, the learned label correlation graph well fits the multilabel classification task though efficiently shimmering the original topological structures amongst labels. Furthermore, the inter-label connections are also inclined by label-specific sample communities (each community for the samples distribution a common label). Namely, if two labels have parallel label-specific sample communities, they are likely to be correlated. Based on this observation, LGL is further regularized by the label HypergraphLaplacian.

**2] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification, in 2011.**
System proposed Classifier chains for multi-label classification itshows that binary relevance-based methods have much to proposition, specifically in positions of scalability to large datasets. System exemplifies this with a novel restraining method that can model label correlations althoughpreserving acceptable computational complexity. Empirical evaluation over a broad range of multi-label datasets with a variety of evaluation metrics exhibits the competitiveness of our chaining techniquein contrast to

related and state-of-the-art methods, together in positions of predictive performance and time complexity.

Based on the binary relevance method, which system argued has many advantages over more sophisticated current methods, specifically in positions of time costs. By passing label correlation information along a chain of classifiers, our technique counteracts the drawbacks of the binary method thoughpreserving acceptable computational complexity. Ancollaborative of classifier chains can be rummage-sale to additional augment predictive performance. Using a range of multi-label datasets and evaluation measures, system accepted out empirical evaluations against a range of algorithms. Our classifier chains technique proved superior to related methods, and in an ensemble scenario was able to improve on state-of-the-art methods, particularly on large datasets. Despite other methods using more complex processes to model label correlations,

### 3] M.-L. Zhang and Z.-H.Zhou, "Multilabel neural networks with applications to functional genomics and text categorization," in 2006.

Proposed Multilabel neural networks with applications to functional genomics and text categorization.It is consequent from the popular Back propagation algorithm through engaging a novel error function capturing the characteristics of multi-label learning, i.e. the labels be appropriate to an case should be graded higher than those not fit in to that instance. Applications to two real world multi-label learning problems, i.e. functional genomics and text categorization, demonstration that the performance of BP-MLL is greater to those of some well-established multi-label learning algorithms**.**

### 4] G. Tsoumakas, I. Katakis, and L. Vlahavas, "Random k-labelsets for multilabel classification in 2011.

System proposed Random k-labelsets for multilabel classification,. System proposed a humble yet in effect multi-label learning method, called label powerset (LP), deliberateseveryseparateamalgamation of labels that be present in the training set as a diverse class value of a single-label classification task. The computational efficiency and predictive performance of LP is tested by application domains with great number of labels and training examples. In these circumstances the number of classes may become very large and all together many classes are associated with very few training examples. To deal with these difficulties, this system recommends breaking the preliminary set of labels into a number of small random subsets, called labelsetsand engaging LP to train a consistent classifier. The label sets can be whichever disjoint or overlapping dependent on which of two approaches is rummage-sale to construct them. The proposed method is called RAkEL (RAndomk lab ELsets), where k is a parameter that requires the size of the subsets. Empirical evidence designates that RAkEL manages to increasesignificantly over LP, specifically in domains with large number of labels and shows competitive performance counter to other high-performing multi-label learning methods. RAkEL could be additional generally thought of as a new approach for creating ancollective of multi-label classifiers by manipulating the label space using randomization. In this sense, RAkEL could be independent of the underlying method for multi-label learning, which in this

.

system is LP. However, system should note that only multi-label learning methods that strongly depend on the specific set of label. Extracting significantsubgraph features, by means ofspecific predefined criteria, to signify a graph in a vectorial space develops a popular solution for graph classification. The most mutualsubgraph selection standard is frequency, which aims to select frequently actingsubgraphs by using frequent subgraph mining methods. For example, one of the maximumgeneral algorithms for frequent subgraph mining is gSpan. Its uses depth first search (DFS) to search most frequent subgraph.

### 5] L Jiao and L Feng proposed Ant Colony optimization in IEEE 2010

Application of an ant colony algorithm for text classification, daily, the amount of information obtainable to us rises. This info would be unusable and not relevant if our ability to efficiently access didn't increase ~~ensembles for classifier benefit,~~ can achieve bet system need tools that permit us to search, sort, index, store, and analyze the available data. For greater benefit, system need tools that permit us to search, sort, index, store, and analyze the available data. System also need tools which assistances us to find desired information in a reasonable time by performing certain tasks for us. One of the promising areas is the automated text classification. Just imagine system have substantial number of texts, which are more simply accessible if they are prepared into typesconferring to their theme. Of course, system can ask a human to categorize the texts by reading them manually. This task is very tough if system do it for hundreds, even thousands of texts. So, it seems necessary to have an automatic text classification application. In this system author presents his experimentations in automated text categorization, where author suggest the use of an ant colony algorithm.

### 6] SasankaPotluri, Christian Diedrich proposed Accelerated Deep Neural Networks for Enhanced Intrusion Detection System IEEE 2016

Basically main focus of this system is to appraise the presentation of the Deep Neural Network (DNN) training related to different processor types and numbers of cores. The acceleration of the training process using the multi core CPU's was quicker than the serial training mechanism. But the GPU's were incapable to achieve the expected performance due to the type of data system used. The NN has use in this scheme for IDS, but most important the system like as multilable verification on each layer for better accuracy purpose.

## 3. PROPOSED SYSTEM

Fig.1 shows the proposed multi-view-graph learning for graph bag classification. In the recommended research work order to further exploit the label correlation information to achieve a cost sensitive classification approach, the edge weights of the label graph are used to weight the slack variables for the relevant-irrelevant label pairs within the ranking NN framework as formulated . Due to the fact that the ranking NN losses in turn provide some structural constraints on the inter-label interactions used for the label graph learning, the label correlation learning task and the classification task are correlated and mutually reinforced. A joint learning of the two tasks should be built to learn the correlation matrix adaptively with the multi label classification problem
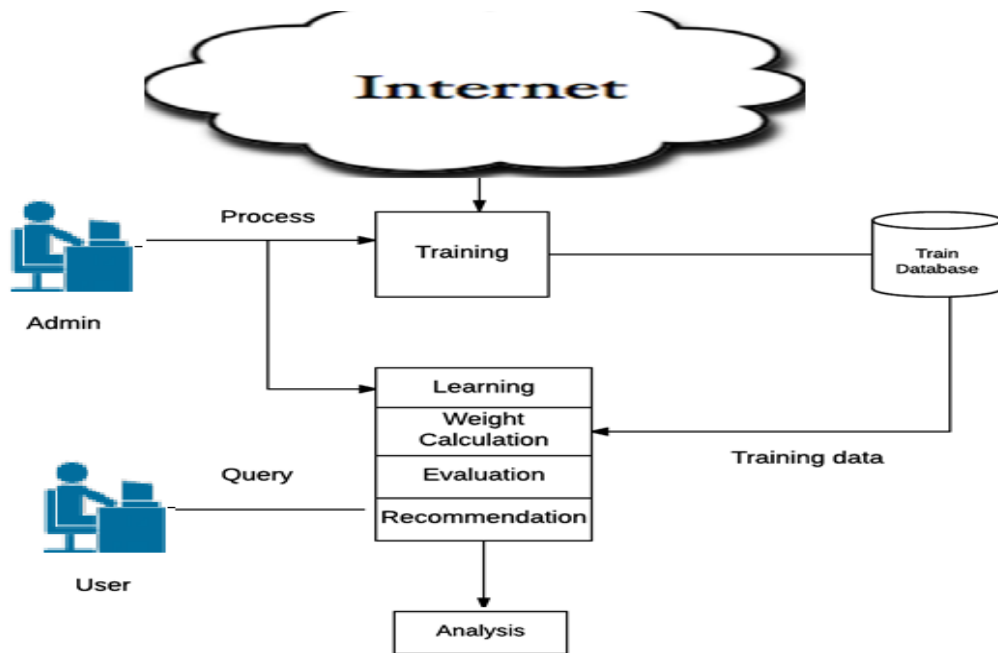
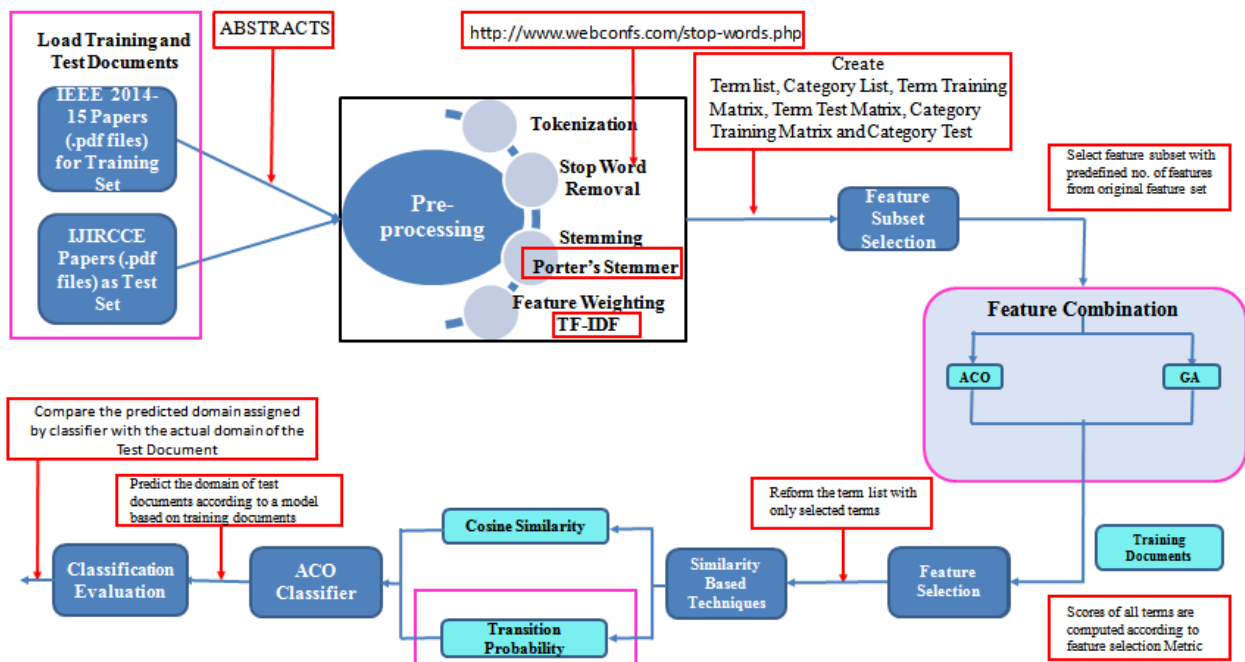**Fig. 1. Proposed system architecture**



**Fig.2: System flow**

In the project we are by standard IEEE 2015 and online image dataset for training as well testing purpose. For the training and testing criteria has been given 70/30 (Training/Testing). Below modules detail shows linear execution of system.

## 3.1 Data Training phase with pre-processing

- This module performs data pre-processing to create train dataset.

- Then first upload the training directory of .pdf dataset and image from NUSWIDE dataset.

- Once upload it will read the data from PDF using PDFBOX API.

- Then tokenization, stop word removal and porter's stemmer will execute.

- Finally TF-IDF will provide the obtainability of current vector and store into feature database

## 3.2 Testing phase with preprocessing and TF-IDF

- First upload the test directory of pdf as well image dataset.

- The initial phase of testing is same like training phase till IDF score calculation.

- Then features are mined using NN

- And classification is done using similarity vector

## 3.3 Feature Selection phase

- This module extract the feature form all buckets using Optimization approach.

- Initial pheromone need to set.

- The pheromone will select the neighbors and strong node for selection.

## 3.4 NN Classification module

- Here NN use for classification purpose.

- Here we find the training dataset with domain detail and feature details.

- Once NN execute it will ask for variation as well generation, after that crossover and mutation execute.

- Finally similarity score will classify each bucket into the respective domain.

## 3.5  Algorithm

### 3.5.1 Stop word Removal Approach

Input: Stop words list L[], String Data D for remove the stop words.

Output: Verified data D with removal all stop words.

Step 1: Initialize the data string S[].

Step 2: initialize a=0,k=0

Step 3:  for each(read a to L)

    If(a.equals(L[i]))

Then Remove S[k]

End for

Step 4: add S to D.

Step 5: End Procedure

### 3.5.2 Stemming Algorithm.

**Input : Word w**

**Output : w with removing past participles as well.**

Step 1: Initialize w

Step 2:  Intialize all steps of Porter stemmer

Step 3: for each (Char ch from w)

    If(ch.count==w.length()) && (ch.equals(e))

    Remove ch from(w)

Step 4: if(ch.endswith(ed))

  Remove 'ed' from(w)

Step 5: k=w.length()

    If(k (char) to k-3 .equals(tion))

   Replace w with te.

Step 6: end procedure

### 3.5.3 TF-IDF

Comment = {c1, c2, c3….cn}
Aspects available in each comment
D = {cmt1, cmt2, cmt3, cmtn}
And comments available in each document
Calculate the Tf score as
tf (t,d) = (t,d)
t=specific term
d= specific document in a term is to be found.

This is known as weight of tf formula for specific comment.

### 3.5.4 Weight calculation Algorithm (NN)

**Input: Each query result table from crawler with CS score, Threshold T for calculate relevancy.**

**Output: classified each attribute with NN classifier with relevancy factor.**

Here we have to find similarity of two vectors: $\vec{a} = (a_1, a_2, a_3, \ldots)$

and $\vec{b} = (b_1, b_2, b_3, \ldots)$, where $a_n$ and $b_n$ are the components of the vector (features of the document, or values for each word of the comment ) and the $n$ is the dimension of the vectors:

$$\vec{a} \cdot \vec{b} = \sum_{i=1}^{n} a_i b_i = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n$$

Step 1: Read each row R from dataset D

Step 2:  for each ( Column c from R)

Step 3:  Get C[i] as category and C[i+1] score

Step 4: summarize all attribute score with sumscore(C)

Step 5: calculate relevancy score for each attribute list.

Step 6: assign each Row class label as relevant as well as irrelevant.

Step 7: Categorize all instances

Step 8: end for end procedure

## Mathematical Model

Let S, be the proposed system which can be represented as

S = {{I}, {$I_t$, $I_s$, $I_{st}$, $I_{tw}$, $I_{tr}$, $I_{ts}$, $I_{fs}$, $I_{cs}$, $I_{tp}$, $I_r$}, {R}}

Where,

I -> Input document as well image collection (for Training and Testing)

$I_t$ -> Abstract reading from input document

$I_s$ ->Applying stop word removal on abstracts

$I_{st}$-> Applying Porter's Stemming on abstract

$I_{tw}$-> Feature Term Set with TF-IDF score

$I_{tr}$-> Training Feature Set

$I_{ts}$-> Test Feature Set

$I_{fs}$ -> Test Feature subset depending upon the training feature set

$I_{cs}$-> Test feature subset with Cosine Similarity

$I_{tp}$-> Test feature subset with Transition Probability

R -> Test document labeled with the appropriate domain/ category

**P** is Learning Algorithm

**Input** = {Text Documents}

**Output**= {Categorized text with their labels}

Where, P represented as **Functions** like Tokenization, Stemming, Stop word Removal, Feature Selection and Feature Transformation.

**P** = {Fx | Input Output}

**Fx** is a function which takes input as text documents and give output as text indexing.

Let S, be the proposed system which can be represented as

$$S = \{\{I\}, \{P\}, \{O\}\}$$

Where,
*I* ->Input document collection (for Training and Testing)
*P* -> Functions used
*O* -> Test document labeled with the appropriate Domain
Where,

$$P = \{f1, f2, f3, f4\}$$

*f1* ->Term Weighting (TF-IDF)
*f2* -> Feature Selection Method (Evaluation Algorithm NN)
*f3* -> Similarity Based Methods (Cosine Similarity and Transition Probability)
*f4* -> Evaluation Parameters (Precision, Recall and Accuracy)

## 4. RESULTS AND DISCUSSION

The below analysis is the system classification graph. The graphs display how system classify the overall inputs into categories. The proposed system is implemented with ACO-GA combination, which gives all results with satisfactory level. For performance assessment 112 documents given for training and 30 documents given for testing. Here system compares the proposed results with two different existing systems.

For the system performance evaluation, calculate the matrices for accuracy. The system is implemented on java 3-tier architecture framework with INTEL 2.8 GHz i3 processor and 4 GB RAM with open source environment. The experimental results illustrate the benefits of multi-view learning methods compared to traditional single-view learning, Table 1 presents a list drawn from several published multi-view learning papers.(Varma and Babu, 2009 ,Zhu et al., 2012 and Rakotomamonjy et al., 2008 used the WebKB data as one of the evaluation datasets

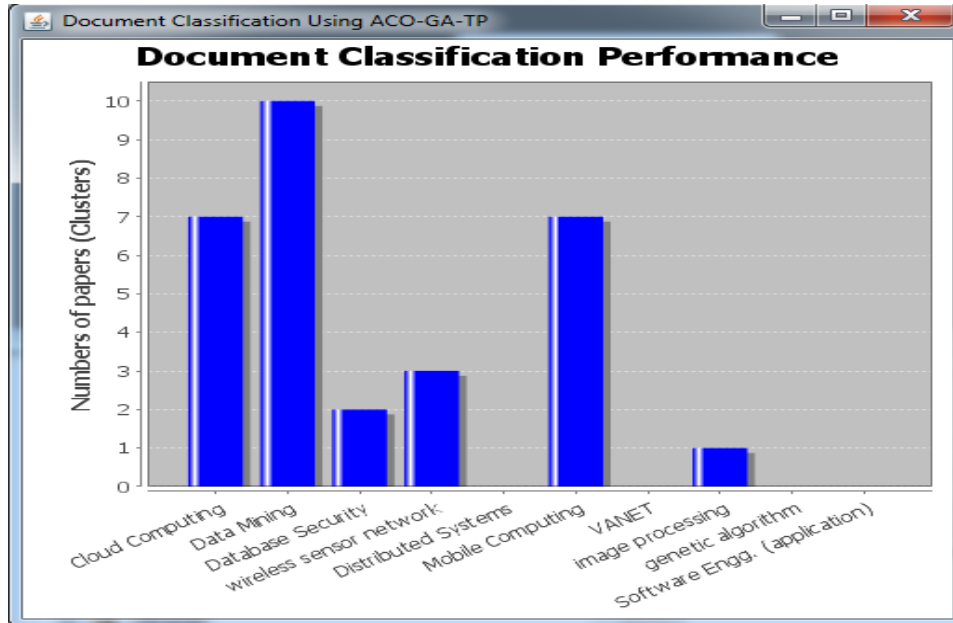| Author | Dataset | Accuracy in % | False Rate in % |
|---|---|---|---|
| Zhu et al., 2012 | Page+Hyperlink | 81.99 | 18.11 |
| Rakotomamonjy et al., 2008 | Wpbc Sonar | 78.50 | 21.50 |
| (Varma and Babu, 2009 | Parkinsons Ionosphere Wpbc | 86.15 | 13.85 |
| Proposed System | IEEE base pdf dataset and NUSWODE dataset (Estimated) | 92.50 | 7.50 |

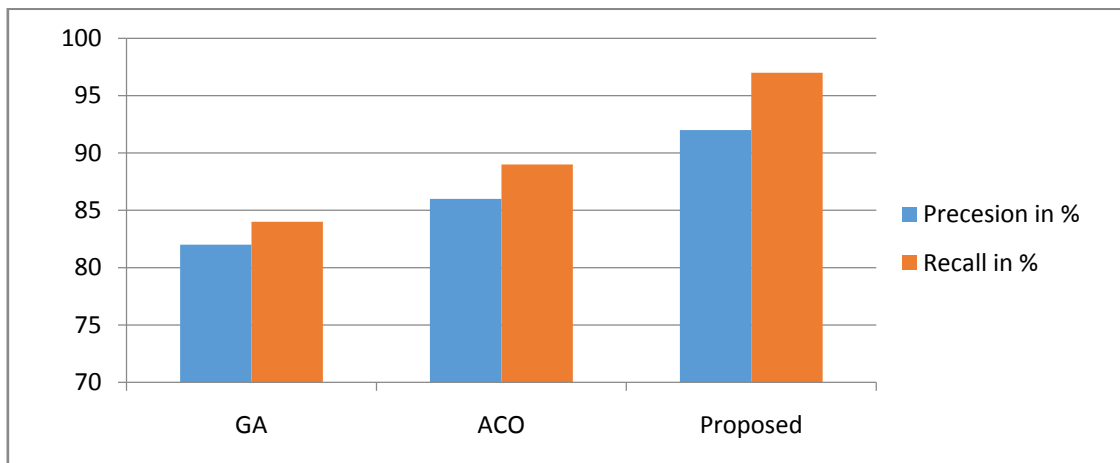**Figure 3: Domain Classification Accuracy**



**Figure 4: System performance results with existing system**

## 5. CONCLUSION

In this work proposed a joint learning scheme for instantaneously modeling label graph learning and multi-label classification. The proposed learning scheme obviously models the inter-label correlations by label graph learning, which is jointly optimized with multi-label classification. As a outcome, the learned label correlation graph is capable of well-fitting the multi-label classification task although efficiently reflecting the underlying topological structures amongst labels. In addition, we have presented a community-aware regularize to capture the context-dependent inter-label interaction information. The proposed work can classify the strong label with test occurrence using NN weight calculation as well classification approach. . Experimental results have verified the efficiency of our approach over several benchmark datasets.

To enhance the system, we have proposed the improved feature selection method by combining ACO and GA algorithms to identify the best minimal feature subset for classification which results in improved accuracy with less computation time. ACO based document classification can be obtained using Cosine Similarity but use of Transition Probability with existing Cosine Similarity improves the accuracy of document classification in terms of Precision and Recall. The proposed system results in a balanced performance of text document classification's accuracy in terms of precision and recall.

## 6. FUTURE WORK

We need to focus for future enhancement for this system

- Sometime system having issue of false result.

- System execution complexity issue when we work with high dimensional or big data.

- System can be work with HDFS framework

# 7. REFERENCES

[1] X. Kong, M. K. Ng, and Z.-H. Zhou, "Transductivemultilabel learning via label set propagation," IEEE Trans. Knowl. Data Eng., vol. 25, no. 3, pp. 704–719, Mar. 2013.

[2] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier chains for multi-label classification," J. Mach. Learn., vol. 85, no. 3, pp. 333–359, Dec. 2011.

[3] M.-L. Zhang and Z.-H.Zhou, "Multilabel neural networks with applications to functional genomics and text categorization," IEEE Trans. Knowl. Data Eng., vol. 18, no. 10, pp. 1338–1351, Oct. 2006.

[4] G. Tsoumakas, I. Katakis, and L. Vlahavas, "Random k-label sets for multilabel classification," IEEE Trans. Knowl. Data Eng., vol. 23, no. 7, pp. 1079–1089, Jul. 2011.

[5] Ant Colony optimization L Jiao, L Feng - Information and Computing (ICIC), 2010 - ieeexplore.ieee.org

[6] A Survey on Approaches of Multirelational Classification Based On Relational Database ShraddhaModi, AmitThakkar, AmitGanatra, International Journal of Engineering and Advanced Technology (IJEAT) ISSN: 2249 – 8958, Volume-1, Issue-3

[7] Lachetar, N. ;Comput. Sci. Dept., Univ. 20 Aout 1955, Skikda, Algeria ; Bahi, H. Application of an ant colony algorithm for text indexing, :Multimedia Computing and Systems (ICMCS), 2011 International Conference –IEEE 2011.