# Multi Label Document classification approach using Machine Learning Techniques: A Survey

## Priyanka B. Sonawane[1],Prof. Pramila M Chawan[2]

*[1]M.Tech Student, Dept of Computer Engineering and IT, VJTI College, Mumbai, Maharashtra, India*
*[2] Associate Professor, Dept of Computer Engineering and IT, VJTI College, Mumbai, Maharashtra, India*

---***---

**Abstract:** - *Multi-label classification is an important machine learning task wherein one assigns a subset of candidate labels to an object. In this paper, we propose a new multi-label classification method based on Conditional Bernoulli Mixtures. Our proposed method has several attractive properties: it captures label dependencies; it reduces the multi-label problem to several standard binary and multi-class problems; it subsumes the classic independent binary prediction and power-set subset prediction methods as special cases; and it exhibits accuracy and/or computational complexity advantages over existing approaches. We demonstrate two implementations of our method using logistic regressions and gradient boosted trees, together with a simple training procedure based on Expectation Maximization. We further derive an efficient prediction procedure based on dynamic programming, thus avoiding the cost of examining an exponential number of potential label subsets. For the testing we will use and show the effectiveness of the proposed method against competitive alternatives on benchmark datasets with pdf. An increasing number of data mining tasks includes the analysis of complex and structured types of data and make use of expressive pattern languages. Most from these applications can't be solved using traditional data mining algorithms. This is the cause of main motivation for the multi-disciplinary field of Multi-Relational Data Mining (MRDM).*

***Key Words*: Text categorization, Machine learning, Multi-Relational Data Mining**

## 1. INTRODUCTION

Text classification aims to classify the text to some predefined categories using a kind of classification algorithm which is performed based on text content. The standard classification corpus has been established and a unified evaluation method is adopted to classify English text based on machine learning which has made a large progress now. In the past ten years management of document based contents (collectively known as information retrieval – IR) have become very popular in the information systems field, due to the greater availability of documents in digital form and resulting need to access them in flexible and efficient ways. Text categorization (TC), the activity of labeling natural language texts with one or more categories from a predefined set, is one such task. Machine learning (ML) methodology, according to which we can automatically build an automatic text classifier by learning, from a set of pre-classified text documents based on the characteristics of the categories of interest. The advantages of this approach is accuracy as compared to that obtained by human beings, and

a considerable savings in terms of expert manpower, since no contribution from either domain experts or knowledge engineers is needed for the construction of the classifier.

Recent years have experienced an increasing number of applications where instances (or samples) are no longer represented by a flat table with instance-feature format, but share some complex structural dependency relationships. Typical examples include XML web pages (i.e. an instance) which point to (or are pointed to) several other XML web pages , a scientific publication with a number of references , posts and responses generated from social networking sites, and chemical compounds with molecules (i.e. nodes) linked together through some bounds (edge) . Given a collection of graph data $\{G_i, y_i\}$, $y_i \in \bullet \{+1,-1\}$, each of which is labeled (+1 for +ve graphs, -1 for -ve graphs), we propose here graph classification model which aims to build prediction model with maximum accuracy in classifying previously unseen graphs. In order to learn classification models, we propose a multi-view-graph learning algorithm (MVGL), which aims to explore sub graph features from multiple graph views for learning, which aims from a set of labeled bags to learn a classifier each having many number of graphs inside the bag. A bag is considered positive, if at least one graph in the bag is positive, and negative otherwise.

Such a multi-graph presentation can be used for large number of real-world applications, such as classification of webpage, where a webpage can be considered as a bag with texts inside it being represented as graphs. Another application is scientific publication classification, where a paper and its references can be represented as a bag of graphs and each graph (i.e., a paper) is formed by using the correlations between keywords in the paper, as shown in Fig. 1. A bag is labelled +ve, if the paper is relevant to a specific topic. Similarly, for online review based product recommendation, each product receives many customer reviews. For each review composed of detailed text descriptions, we can use a graph to represent the review descriptions. Thus, a product can be represented as a bag of graphs. A bag (product) can be labelled as positive if it receives at least one positive review else negative. As a result, we can use graph based classification learning to help recommend products to customers.

### 1.1 Text categorization

Text categorization is the duty of assigning a Boolean value to each pair $<d_i, c_j> \in D \times C$, where D is a domain of documents & $C = \{c_1, . . . , c_{|C|}\}$ is a set of predefined

categories. A value of T assigned to <di ,cj> indicates a decision to file di under cj, while a value of F indicates a decision not to file di under cj. More formally, the task is to approximate the unknown target function ϕ^ : D ×C → {T, F} (that describes how documents ought to be classified) by means of a function ϕ: D × C → {T, F} called the classifier such that ϕ ^ and ϕ "coincide as much as possible".

### 1.2 Single-label vs. multi-label text categorization

Various constraints may be applied on the Text Classification task, depending on the application. For instance we may need that, for a given integer n, exactly n (or ≤ n, or ≥ n) elements of C be assigned to each di∈ D. The case in which exactly 1 category must be assigned to each di ∈ D is often called the case of single-label, while the case in which any number of categories from 0 to |C| may be assigned to the same di ∈ D is dubbed the multi-label case.

### 2. LITERATURE SURVEY

According to [1] proposed Transductive multi label learning via label set propagation, The issue of multilabel characterization has pulled in incredible enthusiasm for the most recent decade, where every case can be relegated with an arrangement of various class marks at the same time. It has a wide assortment of true applications, e.g., programmed picture explanations and quality capacity examination. Ebb and flow research on multilabel arrangement concentrates on administered settings which expect presence of a lot of named preparing information. Be that as it may, in numerous applications, the marking of multileveled information is greatly costly and tedious, while there are frequently rich unlabeled information accessible. This paper, they examine the issue of transductive multilabel learning and propose a novel arrangement, called TRAsductive Multilabel Classification (TRAM), to successfully allot an arrangement of numerous names to every occasion. Not the same as administered multilabel learning techniques, system evaluate the mark sets of the unlabeled cases successfully by using the data from both marked and unlabeled information. System first plans the transductive multilabel learning as an enhancement issue of evaluating name idea pieces. At that point, it infer a shut structure answer for this improvement issue and propose a compelling calculation to dole out name sets to the unlabeled examples. Observational studies on a few certifiable multilabel learning assignments exhibit that our TRAM strategy can successfully support the execution of multilabel order by utilizing both marked and unlabeled information.

System first, formulates the task as an optimization problem which is able to exploit unlabeled data to obtain an effective model for assigning appropriate multiple labels to instances. Then, develop an efficient algorithm which has a closed-form solution for this optimization problem. Empirical studies on a broad range of real-world tasks demonstrate that our TRAM method can effectively boost the performance of multilabel classification by using unlabeled data in addition to labeled data.

According to [2] proposed Classifier chains for multi-label classification it shows that binary relevance-based methods have much to offer, especially in terms of scalability to large datasets. System exemplify this with a novel chaining method that can model label correlations while maintaining acceptable computational complexity. Empirical evaluation over a broad range of multi-label datasets with a variety of evaluation metrics demonstrates the competitiveness of our chaining method against related and state-of-the-art methods, both in terms of predictive performance and time complexity.

Based on the binary relevance method, which system argued has many advantages over more sophisticated current methods, especially in terms of time costs. By passing label correlation information along a chain of classifiers, our method counteracts the disadvantages of the binary method while maintaining acceptable computational complexity. An ensemble of classifier chains can be used to further augment predictive performance. Using a variety of multi-label datasets and evaluation measures, we carried out empirical evaluations against a range of algorithms. Our classifier chains method 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,    ensembles of classifier chains can achieve better predictive performance and are efficient enough to scale up to very large problems.

According to [3], proposed Multilabel neural networks with applications to functional genomics and text categorization. It is derived from the popular Back propagation algorithm through employing a novel error function capturing the characteristics of multi-label learning, i.e. the labels belonging to an instance should be ranked higher than those not belonging to that instance. Applications to two real world multi-label learning problems, i.e. functional genomics and text categorization, show that the performance of BP-MLL is superior to those of some well-established multi-label learning algorithms**.**

System [4] proposed Random k-label sets for multilabel classification,. System proposed a simple yet effective multi-label learning method, called label power set (LP), considers each distinct combination of labels that exist in the training set as a different class value of a single-label classification task. The computational efficiency and predictive performance of LP is challenged by application domains with large number of labels and training examples. In these cases the number of classes may become very large and at the same time many classes are associated with very few training examples. To deal with these problems, this paper proposes breaking the initial set of labels into a number of small random subsets, called label sets and employing LP to train a corresponding classifier. The label sets can be either disjoint or overlapping depending on which of two strategies is used to construct them. The proposed method is called RA*k*EL (RAndom *k* labELsets), where *k* is a parameter that specifies the size of the subsets. Empirical evidence indicate

that RA*k*EL manages to improve substantially over LP, especially in domains with large number of labels and exhibits competitive performance against other high-performing multi-label learning methods. RA*k*EL could be more generally thought of as a new approach for creating an ensemble of multi-label classifiers by manipulating the label space using randomization. In this sense, RA*k*EL could be independent of the underlying method for multi-label learning, which in this paper is LP. However, we should note that only multi-label learning methods that strongly depend on the specific set of label.

System [5] "gSpan: Graph-based substructure pattern mining,". Extracting important subgraph features, using some predefined criteria, to represent a graph in a vectorial space becomes a popular solution for graph classification. The most common subgraph selection criterion is frequency, which intends to select frequently appearing subgraphs by using frequent subgraph mining methods. For example, one of the most popular algorithms for frequent subgraph mining is gSpan [5]. Its uses depth first search (DFS) to search most frequent subgraph.

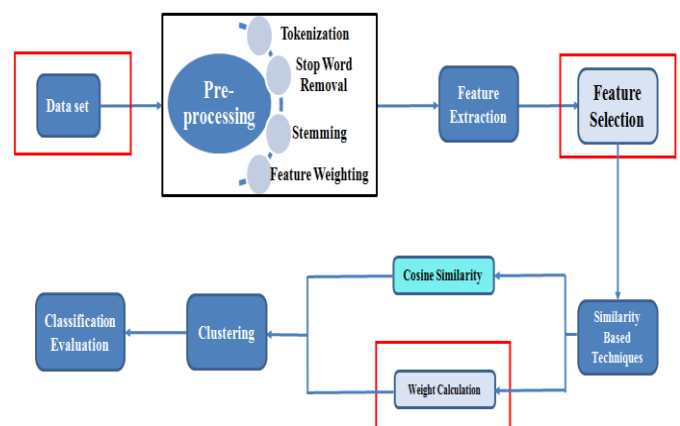| id | Title | Method Used | Gaps Identification |
|---|---|---|---|
| 1 | Transductive multilabel learning via label set propagation | Rule and signature base supervised learning | Single instance can be classify multiple classes |
| 2 | Classifier chains for multi-label classification | binary relevance-based methods used the term scalability | Generate very high computational complexity |
| 3 | Multilabel neural networks with applications to functional genomics and text categorization | Used neural network algorithm named BP-MLL, i.e., back propagation for multilabel learning | High false positive ratio. |
| 4 | Random k-label sets for multi label classification | effective multilabel learning method, called label power set (LP) and RAkEL (RAndom k labELsets) | It can support structured dataset only |
| 5 | gSpan: Graph-based substructure pattern mining | gSpan Graph base structured Pattern Mining | It loss some useful information during the graph construction |
| 6 | Multilabel classification with label correlations and missing labels | multilabel classification with label correlations and missing labels" | Much low accuracy when data contains the missing labels |
| 7 | Multi-label learning by exploiting label dependency | Bayesian network structure | Accuracy and efficiency issue |
| 8 | The landmark selection method for multiple output prediction | Landmark selection method | Pattern values missing during the classification |
| 9 | Multi-label learning by exploiting label correlations locally | ML-LOC : Multi-Label learning using Local Correlation | Time complexity issues during the image features extraction |
| 10 | Efficient multi-label classification with many labels | Multi-label classification via CSSP (ML-CSSP). | Error rate is very high during the testing |

## 3 PROPOSED SYSTEM



**Figure 1. Proposed system architecture**

The proposed system works with multi label classification approach using machine learning algorithm. The system works based on the semi supervised learning. In raining phase system first create the Background Knowledge (BK) and testing phase works for classification using specific weight of desired class label.

## 4. CONCLUSION

In the proposed survey we evaluate the various methods of classification. The multi label classification issue has occurred in many existing systems, also some systems having time complexity issues. The proposed work can classify the strong label with test instance using machine learning weight calculation as well classification approach. The proposed learning scheme explicitly models the inter-label correlations by label graph learning, which is jointly optimized with multilabel classification. As a result, the learned label correlation graph is capable of well fitting the multilabel classification task while effectively reflecting the underlying topological structures among labels.

## REFERENCES

[1] X. Kong, M. K. Ng, and Z.-H. Zhou, "Transductive multilabel 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-labelsets for multilabel classification," IEEE Trans. Knowl. Data Eng., vol. 23, no. 7, pp. 1079–1089, Jul. 2011.

[5] X. Yan and J. Han, "gSpan: Graph-based substructure pattern mining," in Proc. 2nd ICDM, Washington, DC, USA, 2002, pp. 721–724.

[6] W. Bi and J. T. Kwok, "Multilabel classification with label correlations and missing labels," in Proc. Assoc. Adv. Artif. Intell., 2014, pp. 1680– 1686.

[7] M.-L. Zhang and K. Zhang, "Multi-label learning by exploiting label dependency," in Proc. ACM SIGKDD Conf., 2010, pp. 999–1008.

[8] K. Balasubramanian and G. Lebanon, "The landmark selection method for multiple output prediction," in Proc. Int. Conf. Mach. Learn., 2012, pp. 983–990.

[9] S.-J. Huang and Z.-H. Zhou, "Multi-label learning by exploiting label correlations locally," in Proc. Assoc. Adv. Artif. Intell., 2012, pp. 949– 955.

[10] W. Bi and J. Kwok, "Efficient multi-label classification with many labels," in Proc. Int. Conf. Mach. Learn., 2013, pp. 405–413.