# Addressing Multilabel Classification Problem Via Co-evolutionary Learning Algorithm

**Gayatri T. Urade[1], Prof. Pravin G. Kulurkar[2]**

M.Tech CSE, Vidarbha Institute of Engineering, Nagpur[1]

H.O.D, CSE, Vidarbha Institute of Engineering, Nagpur[2]

**Abstract:** Multi-label classification refers to the task of predicting potentially multiple labels for a given instance. Conventional multi-label classification approaches focus on the single objective setting, where the learning algorithm optimizes over a single performance criterion (e.g.Ranking Loss) or a heuristic function. The basic assumption is that the optimization over one single objective can improve the overall performance of multi-label classification and meet the requirements of various applications. However, in many real applications, an optimal multi-label classifier may need to consider the tradeos among multiple conflicting objectives, such as minimizing Hamming Loss and maximizing Micro F1.

**Keywords:** Muti Label, Web, Learning, Memory efficiency.

## 1. INTRODUCTION

Data classification is one of the major issues in data mining and machine learning. Generally speaking, it consists of two stages, that is, building classification models and predicting labels for unknown data. Depending on the number of labels tagged on each data, the classification problems can be divided into single-label and multilabel classification.

In the former, the class labels are mutually exclusive and each instance is tagged with only one class label. On the contrary, each instance may be tagged with more than one class label simultaneously. The multilabel classification problems are ubiquitous in real-world applications, such as text categorization, image annotation, bioinformatics, and information retrieval.

Multilabel learning is now receiving an increasing attention from a variety of domains and many learning algorithms have been witnessed. Similarly, the multilabel learning may also suffer from the problems of high dimensionality, and little attention has been paid to this issue. In this paper, we propose a new ensemble learning algorithms for multilabel data. The main characteristic of our method is that it exploits the features with local discriminative capabilities for each label to serve the purpose of classification. Specifically, for each label, the discriminative capabilities of features on positive and negative data are estimated, and then the top features with the highest capabilities are obtained. Finally, a binary classifier for each label is constructed on the top features. Experimental results on the benchmark data sets show that the proposed method outperforms four popular and previously published multilabel learning algorithms.

Our Moml (multi-objective multi-label classification) algorithm finds a set of non-dominated solutions which are optimal according to the different tradeoffs of the multiple objectives. So users can flexibly construct various combined predictive models from the solution set, which helps to provide more meaningful classification results in different application scenarios. Empirical studies on real-world tasks demonstrate that the Moml (multi-objective multi-label classification) can effectively boost the overall performance of multi-label classification, not limiting to the optimization objectives

Aim is to introduce a technique called Moml, (multi-objective multi-label classification) which is an unsupervised proposal that learns extraction rules from a set of web documents that were generated by the same server-side template. It builds on the hypothesis that shared patterns are not likely to provide any relevant data and are, thus, part of the template. Whenever it finds a shared pattern, it partitions the input documents into the prefixes, separators and suffixes that they induce and analyses the results recursively, until no more shared patterns are found. Prefixes, separators, and suffixes are organised into a Moml (multi-objective multi-label classification) tree that is later traversed to build a regular expression with capturing groups that represents the template that was used to generate the input documents.

Thanks to the capturing groups, the expression can be used to extract data from similar documents. Note that our technique does not require the user to provide any multilabel; instead, he or she must interpret the resulting regular expression and map the capturing groups that represent the information of interest onto the appropriate structures.

## 2. LITERATURE SURVEY

Traditional binary or multi-class classification can be regarded as degenerated version of multilabel classification if each instance is confined to have only one

class label [6]. The generality of multilabel classification makes it challenging to solve. The main challenge of multilabel classification lies in how to effectively and efficiently exploit correlations among labels [7]. A major difference between multilabel classification and traditional binary or multi-class classification is that labels in multilabel classification are not mutually exclusive but may be correlated. The correlations among labels are beneficial to label pre-diction, however, to exploit correlations among labels is a nontrivial task.

Firstly, labels may be correlated in various degrees including low-order and high-order correlations. Secondly, the number of possible label combinations is exponential to the number of labels, which can lead to great challenges to the efficiency and scalability of the learning algorithm when the number of labels becomes large.

In recent years, many learning algorithms have been proposed in the literature to deal with multilabel classification problems [6], [8], [9]. However, most of the existing multilabel learning algorithms cannot be both effective and efficient in exploiting correlations among labels. For example, binary relevance (BR) [8] is conceptually simple but it ignores the correlations among labels, other multilabel learning algorithms such as calibrate label ranking (CLR) [10], random k-labelset (RAkEL) [11] and classifier chains (CC) [12] consider the correlations among labels, however, their computational complexity grows dramatically as the number of labels increases. Thus, how to effectively and efficiently make use of correlations in the label space still remains an open question.
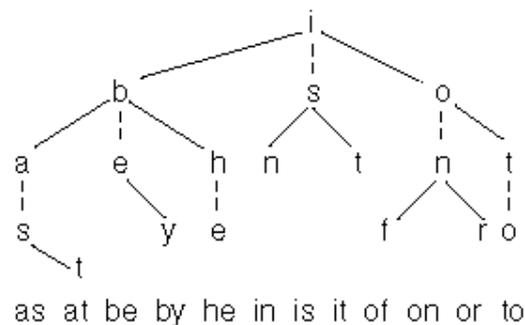
Modeling dependencies among class labels is one of the core goals of multilabel classification [17], [18]. Zhang and Zhou [6], [19] categorize the existing multilabel classification algorithms into three families, i.e. first-order strategy, second-order strategy, and high-order strategy based on the order of label correlations that the learning algorithms have considered. Recently, several dimensionality reduction algorithms have been proposed for multilabel learning [7], [20]. Many of these algorithms attempt to solve the curse of dimensionality problem meanwhile capture correlations among labels.

The first-order strategy treats class labels independently and ignores correlations among labels. Binary relevance (BR) [8] and multilabel k nearest neighbors (ML-kNN) [21] are two representative algorithms of the first-order strategy. BR deals with multilabel learning problem by decomposing it into a lot of binary learning problems (one for each label) and learning each binary classifier independently. Because of its conceptual simplicity, BR is usually used as building blocks in many advanced multilabel learning methods. The basic idea of ML-kNN is to adapt k nearest neighbor techniques to deal with multilabel data, where maximum a posterior (MAP) rule is utilized to make pre-diction by reasoning with the labeling information embo-died in the neighbors. The second-order

strategy exploits the correlations between any two different labels

## 3. PROPOSED METHODOLGY

It's a tree data structure in which each node can have three child nodes.the nodes are divide as left, mid, right. comparing to other tree,such us binary tree .in that all nodes on left child have smaller value compare to the right child. the Moml search tree low and high poiners are shown angles line ,while equal pointers are shown as vertical lines. To calculate the maximum nodes they uses,

If a node occupies TREE ,[K]
Left Child is stored in TREE [3K-1]
Mid Child is stored in TREE .[3K]
Right Child is stored in TREE .[3K+1]



### A. Random Process
A random ―process checking greatly reduces the workload of services. Thus, a probabilistic automatic on sampling checking is preferable to realize the secret key manner, as well as to rationally allocate resources and non repeat keywords.

An efficient algorithm is used to since the single sampling checking may overlook a small number of data abnormalities.

### B. Unsupervised data Classification
Unsupervised data classification refers to the pages in a web site so that each cluster includes a set of web pages that can be classified using a unique class. We propose CALA, a new automated proposal to generate URL-based web page classifiers.
Our proposal builds a number of URL patterns that represent the different classes of pages in a web site, so further pages can be classified by matching their URLs to the patterns.

### C. Capturing Groups
It finds a shared pattern, it partitions the input documents into the prefixes, separators and suffixes that they induce and analyses the results recursively, until no more shared patterns are found. Prefixes, separators, and suffixes are organized into a Moml tree that is later traversed to build a regular expression with capturing groups that represents

the template that was used to generate the input documents. Thanks to the capturing groups, the expression can be used to extract data from similar documents.

### D. Relevant data

The wrapper generation in this type of data set is more challenging since there is no inherent measurement of data mining for discovering rare events.

The relevant data is especially challenging because of the difficulty of defining a data for categorical data or combination of relevant and irrelevant data. Automatic wrapper generation can be implemented as a preprocessing step prior to the application of an identifying the relevant data.

### E. Web data extractors

Web data extractors are used to extract data from web documents in order to feed automated processes. In this article, we propose a technique that works on two or more web documents generated by the same server side template and learns a regular expression that models it and can later be used to extract data from similar documents.

Web data extractors to work on proposals to learn those automatically using supervised techniques. We have presented an effective and efficient unsupervised data extractor called Trinity.

It is based on the hypothesist that web documents generated by the same server-side template share patterns that do not provide any relevant data, but help delimit them. The rule learning algorithms arches for these patterns and creates a Moml tree, which is then used to learn regular expression that represents the template that was used to generate input web documents. using the Euclidean distance which finds the minimum difference between the weights of the input image and the set of weights of all images in the database.

## 4. IMPLEMENTATION DETAILS

Three major steps of this approach can be integrated as follows:-

a. Web-log data collection: The logs we research are of W3C Extended Log File Format under KDD Cup dataset. Web log networking or medical data is collected from the server of website for the period of one month for experimental purpose

b. Data pre-processing: We can use database software Access and Java programming language to implement the preprocessing work. Also web-log file preprocessing tools such as WEBMINER, AWStat can be used for data cleaning, user identification and path completion.

c. Web–usage mining from web-log files:
The final step of web-usage mining can be implemented using neural network approach via. Back Propagation algorithm
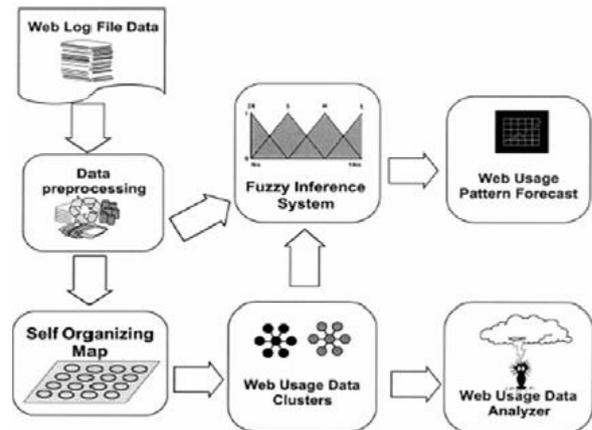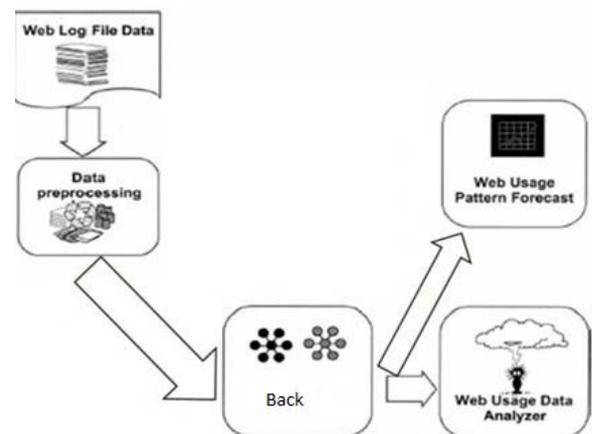


Fig 2: Moml Approach for Web Mining



Figure 3: Reduction of Stages on Moml after Back Propagation implementation.

If any Web-mining researches apply this Moml, then can easily obtained best result than any implemented Web mining techniques because of vigilance parameter, top-down and bottom-up weights.

Also using Moml, it is more beneficiary to minimize the number of steps in Web mining as compare to neuro-fuzzy approach. As neuro-fuzzy approach uses five major steps to produce the Webusage pattern forecast, and Web-usage data analyzer; named Web-log data collection, data preprocessing, self-organizing map, Web-usage data cluster, and fuzzy inference system (Figure. 2). But Moml use only three steps as Web-log data collection, data preprocessing, and Backpropagations itself (Figure 3).

## 5. CONCULSION

In this paper, we propose a Trinary method to deal with multilabel classification problems. We firstly convert the conventional hyper network into a multilabel hyper-network. We then propose a co-evolutionary learning algorithm to learn a multilabel hypernetwork from training data. Experimental results on a variety of multilabel. Data sets show that Trinary is effective then Co-MLHN in addressing multilabel classification problems and scales well with respect to the number of labels.

In Co-MLHN, we only consider the inclusive correla-tions among labels, i.e. the co-occurrence of labels. How-ever, the exclusive relationship among labels is also im-portant to multilabel classification. In the future, we will study how to exploit both inclusive and exclusive correla-tions among labels using hypernetwork.

## REFERENCES

[1] S. Gao, W. Wu, C. H. Lee, and T. S. Chua, "A MFoM Learning Approach to Robust Multiclass Multi-Label Text Categoriza-tion," in Proc. 21th International Conference on Machine Learning, pp. 406-417, 2004.

[2] J. Y. Jiang, S. C. Tsai, and S. J. Lee, "FSKNN: Multi-Label Text Categorization Based on Fuzzy Similarity and k Nearest Neighbors," Expert System with Application, vol. 39, no. 3, pp. 2813-2821, 2012.

[3] M. R. Boutell, J. Luo, X. Shen, and C. M. Brown, "Learning Multilabel Scene Classification," Pattern Recogniti., vol. 37, no. 9, pp. 1757-1771, 2004.

[4] G. Qi, X. Hua, Y. Rui, J. Tang, T. Mei, and H. Zhang, "Correla-tive Multi-Label Video Annotation," in Proc. 15th Int. Conf. Multimedia, pp. 17-26, 2007.

[5] Elisseeff and J. Weston, "A Kernel Method for Multi-Labeled Classification," in Advances in Neural Information Processing Systems 14, T. G. Dietterich, S. Becker and Z. Gha-hramani, Eds. Cambridge, MA, USA: MIT Press, pp. 681-687, 2002.

[6] M. L. Zhang and Z. H. Zhou, "A Review on Multi-Label Learn-ing Algorithms," IEEE Trans. Knowl. Data Eng., vol. 26, no. 8, pp. 1819-1837, 2014.

[7] L. Sun, S. W. Ji, and J. P. Ye, Multi-Label Dimensionality Reduc-tion. CRC Press, pp. 8-9, 2014.

[8] G. Tsoumakas and I. Katakis, "Multi-Label Classification: An Overview," Int. J. Data Warehousing and Mining, vol. 3, no. 3, pp. 1-13, 2007.

[9] G. Madjarov, D. Kocev, D. Gjorgjevikj, and S. Dzeroski, "An Extensive Experimental Comparison of Methods for Multi-Label Learning," Pattern Recogniti., vol. 45, no. 9, pp. 3084-3104, 2012.

[10] J. Fürnkranz, E. Hüllermeier, E. Loza Mencía, and K. Brinker, "Multilabel Classification via Calibrate Label Ranking," Mach. Learn., vol. 73, no. 2, pp. 133-153, 2008.

[11] G. Tsoumakas, I. Katakis, and I. Vlahavas, "Random k-labelset for Multi-Label Classification," IEEE Trans. Knowl. Data Eng., vol. 23, no. 7, pp. 1079-1089, 2011.

[12] J. Read, B. Pfahringer, G. Holmes, and E. Frank, "Classifier Chains for Multi-Label Classification," in Lecture Notes in Artifi-cial Intelligence 5782, W. Buntine, M. Grobelnik, and J. Shawe-Taylor, Eds. Berlin, Germany: Springer, pp. 254-269, 2009.

[13] B. T. Zhang, "Hypernetworks: A Molecular Evolutionary Ar-chitecture for Cognitive Learning and Memory," IEEE Compu-tational Intelligence Magazine, vol. 3, no. 3, pp. 49-63, 2008.

[14] J. K. Kim and B. T. Zhang, "Evolving Hypernetworks for Pat-tern Classification," in Proc. IEEE Congress on Evolutionary Com-putation (CEC 2007), pp. 1856-1862, 2007.

[15] E. S. Kim, J. W. Ha, W. H. Jung, J. H. Jang, J. S. Kwon, and B. T. Zhang, "Mutual Information-Based Evolution of Hypernetwork for Brain Data Analysis," in Proc. IEEE Congress on Evolutionary Computation (CEC 2011), pp. 2721-2727, 2011.