

Discovering Frequent Patterns From Web Log based Unordered Unrooted Tree

Mr. Dhananjay G. Telavekar¹ And Mr.H.A.Tirmare²

^{1,2} Department of Technology And Assistant Professor, Shivaji University Kolhapur

Abstract— Mining frequent trees is very useful in domains like bioinformatics, web mining, mining semi-structured data, etc. The proposed frequent restrictedly embedded sub tree miner (FRESTM), is an efficient algorithm for mining frequent, unordered, embedded sub-trees in a database of labeled trees. The key contributions of our work are as follows: The algorithm enumerates all embedded, unordered trees. A new equivalence class extension scheme generates all candidate trees. The notion of scope-list joins is extended to compute the frequency of unordered trees. The performance evaluation on several synthetic and real world data shows that FRESTM is an efficient algorithm, which has performance comparable to TreeMiner, that mines only ordered trees.

Keywords— Tree Mining, Embedded Trees, Unordered Trees, Pattern Mining.

I. INTRODUCTION

In recent researches the main problem of finding frequent patterns from a database of graphs has several important applications in different areas like bioinformatics, user web log analysis, semi-structured XML data [12], web mining, RNA, phylogeny [15], prerequisite trees, and chemistry compound data [9], and network routing.

A fundamental problem in many other data mining tasks such as association rule mining, classification and clustering. Although finding of frequent patterns (for example closed patterns) has found more interest, developing efficient algorithms for finding frequent patterns is still important, because the efficiency of the algorithms of finding condensed representations depends on the efficiency of the frequent pattern mining algorithms.

Whereas item set mining and sequence mining have been studied extensively in the past, recently there has been tremendous interest in mining increasingly complex pattern types such as trees and graphs. For example several algorithms for tree mining have been proposed which include TreeMiner [5], which mines embedded, ordered trees, FreqT which mines induced ordered trees, FreeTreeMiner which mines induced, unordered, free trees (i.e., there is no distinct root); TreeFinder which mines embedded, unordered trees (but it may miss some patterns; it is not complete); and PathJoin, uFreqt [14], uNot [18], CMTreMiner [4], and Hybrid Tree Miner which mine induced, unordered trees. Our focus in this paper is on a complete and efficient algorithm for mining frequent, labeled, rooted, unordered, unrooted, embedded subtrees and graphs. An efficient algorithm is introduced for the problem of mining frequent, unordered, embedded sub trees in a database of trees. The key contributions of our work are as follows:

- The first algorithm enumerates all embedded, unordered trees.
- A new self contained equivalence class extension scheme generates all candidate trees. Only potentially frequent extensions are considered, but some redundancy is allowed in the candidate generation to make each class self contained.
- The notion of scope-list joins is extended for fast frequency computation for unordered trees. Performance evaluation is conducted on several synthetic dataset and a real web log dataset to show that an efficient algorithm, which has performance comparable to TreeMiner that mines only ordered trees.

II. LITERATURE SURVEY

In this section, the references are collected from all conferences, sites, articles, books from internet which helps to implement the project. For development of this project we referred some of Base papers, Ideas which helps in development, testing, and in deployment phase.

Sen Zhang, Zhihui Du, and Jason T. L. Wang, members of IEEE have proposed a transaction paper [1] on “New Techniques for Mining Frequent Patterns in Unordered Trees” which is the main base paper of project. The paper is all about tree mining problem that aims to discover restrictedly embedded subtree patterns from a set of rooted labeled unordered trees. And the properties of a canonical form of unordered trees, and develop new Apriori-based techniques to generate all candidate subtrees level by level through two efficient rightmost expansion operations: 1) pairwise joining and 2) leg attachment. Also restrictedly embedded subtree detection can be achieved by calculating the restricted edit distance between a candidate subtree and a data tree. These techniques are then integrated into an efficient algorithm, named frequent restrictedly embedded subtree miner (FRESTM), to solve the tree mining problem.

Mostafa Haghiri Chehrehgani, Morteza Haghiri Chehrehgani, Caro Lucas, and Masoud Rahgozar [2] present OInduced, which is a novel and efficient algorithm for finding frequent ordered induced tree patterns from a database of rooted ordered trees. First, log data are converted into rooted ordered trees, and a set of frequent patterns is extracted from them. Then, based on these patterns, a structural classifier is built to classify different users. Structural classifiers show higher performance compared to traditional classifiers, which treat each tree as a bag of words.

Sen Zhang and Jason T.L.Wang [3] puts forth framework for tackling the FAST problem for both rooted and unrooted phylogenetic trees using data mining techniques. We first develop a novel canonical form for rooted trees together with a phylogeny-aware tree expansion scheme for generating candidate subtrees level by level. Then, we present an efficient algorithm to find all FASTs in a given set of rooted trees, through an Apriori-like approach.

Yi Xia, Yirong Yang, Richard R. Muntz, and Yun Chi [4] proposed CMTreeMiner, a computationally efficient algorithm that discovers only closed and maximal frequent sub-trees in a database of labeled rooted trees, where the rooted trees can be either ordered or unordered. The algorithm mines both closed and maximal frequent sub-trees by traversing an enumeration tree that systematically enumerates all frequent sub-trees.

Mohammed J. Zaki [5] present an example for tree mining, consider the problem of mining structural patterns in a data set of Ribonucleic acid (RNA) molecules, which can be represented as trees. To get information about a newly sequenced RNA, researchers may compare it with known RNA structures, looking for common topological patterns, which provide important clues to the function of the RNA.

K. G. Khoo and P. N. Suganthan [6] proposed A genetic algorithm (GA)-based optimization procedure for structural pattern recognition in a model-based recognition system using attributed relational graph (ARG) matching technique. Our work is to improve the GA-based ARG matching procedures leading to a faster convergence rate and better quality mapping between a scene ARG and a set of given model ARGs.

C. H. Leung and Ching Y. Suen [7] A top-down elastic approach to pattern matching and its application to complex handwritten Chinese character recognition is discussed.

Dennis Shasha, Jason Tsong-Li Wang, Kaizhong Zhang and Frank Y. Shih [8] Presents an efficient enumerative algorithm and several heuristics leading to approximate solutions. The algorithms are based on probabilistic hill climbing and bipartite matching techniques.

Jason Tsong-Li Wang, Karpjoo Jeong, and Dennis Shasha [9] presents Approximate-Tree-By-Example (ATBE), which allows inexact matching of trees. The ATBE system interacts with the user through a simple but powerful query language; graphical devices are provided to facilitate inputting the queries.

Yun Chi, Yirong Yang, Richard R. Muntz [10] was proposed Hybrid Tree Miner, a computationally efficient algorithm that discovers all frequently occurring sub-trees in a database of rooted unordered trees. The algorithm mines frequent sub-trees by traversing an enumeration tree that systematically enumerates all sub-trees. The enumeration tree is defined based on a novel canonical form for rooted unordered trees—the breadth-first canonical form (BFCF).

Aída Jiménez, Fernando Berzal, Juan-Carlos Cubero [11] was proposed A scalable and parallelizable algorithm to mine partially-ordered trees. Our algorithm, POTMiner, is able to identify both induced and embedded subtrees in such trees. As special cases, it can also handle both completely ordered and completely unordered trees.

Mong Li Lee, Liang Huai Yang, Wynne Hsu, Xia Yang [12] develop a technique to determine the degree of similarity between DTDs. Our similarity comparison considers not only the linguistic and structural information of DTD elements but also the context of a DTD element (defined by its ancestors and descendants in a DTD tree).

Lizhi Liu, Jun Liu [13] The topic incorporates an efficient algorithm for mining frequent, ordered, embedded subtree in tree-like databases. Using a new data structure called scope-list, which is a canonical representation of tree node, the algorithm first generates all candidate trees, then enumerates embedded, ordered trees, finally joins scope-list to compute frequency of embedded ordered trees.

Siegfried Nijssen and Joost N. Kok [14] edifies on frequent tree mining. The frequent tree discovery task is the task of discovering all trees referred to as the pattern trees that occur frequently in some large tree called a data tree.

Jason T. L. Wang, Dennis Shasha [15] present a new FSM technique for finding patterns in rooted unordered labeled trees. A rooted unordered labeled tree is a tree in which there is a root for the tree, each node may have a label, and the left-to-right order among siblings is unimportant.

Lei Zou, Yansheng Lu, Huaming Zhang [16] propose a problem to discover frequent induced subtree patterns that are super trees of a given pattern tree specified by users, i.e. frequent induced subtree patterns with subtree-constraint. Most existing frequent subtree mining algorithms are based on right-most extension, which does not work well in the new problem. So free extension is presented to replace right-most extension.

Tatsuya Asai, Hiroki Arimura, Takeaki Uno, and Shin-ichi Nakano [17] proposed and presents efficient algorithm for discovering frequent substructures in a large graph structured data and collection of semi-structured data, where both of the patterns and the data are modeled by labeled unordered trees.

III. CHOICE OF TOPIC WITH REASONING

Frequent patterns are itemsets, subsequences, or substructures that appear in a data set with frequency no less than a user-specified threshold. For example, a set of items, such as milk and bread, that appear frequently together in a transaction data set, is a frequent itemset. A subsequence, such as buying first a PC, then a digital camera, and then a memory card, if it occurs frequently in a shopping history database, is a (frequent) sequential pattern.

A substructure can refer to different structural forms, such as subgraphs, subtrees, or sub lattices, which may be combined with itemsets or subsequences. If a substructure occurs frequently in a graph database, it is called a (frequent) structural pattern. Finding frequent patterns plays an essential role in mining associations, correlations, and many other interesting relationships among data. Moreover, it helps in data indexing, classification, clustering, and other data mining tasks as well. Thus, frequent pattern mining has become an important data mining task and a focused theme in data mining research.

IV. OUTLINE OF WORK

The proposed system consists of following modules:

Module 1-Preprocessing and Tree Generation:

In preprocessing web log defining will be done this include removing incomplete web log, reducing noisy data and data set conversion. Tree generation will convert session web logs to tree structure the session web logs are in form of associated manner.

Module 2-Canonical Representation:

An unordered tree t is in its canonical form if no equivalent ordered tree t' exists with $dls(t') < dls(t)$, that is the canonical form of an unordered tree should result in the lightest dls among all of its equivalent ordered trees. Directly removing the last node of a canonicalized tree t will result in a residue tree still in its canonical form. Here directly removing means removing a node without further canonicalizing the resulting tree. Therefore, if t is an unordered tree in its canonical form, then every downward sub-tree and every prefix of t is also automatically in its canonical form.

Module 3-Support Counting:

To count the support, calculate the occurrence number, of a candidate k -subtree pattern in the whole data set, intuitively, we should run the restrictedly embedding detection subroutine on the candidate pattern tree against all data trees one by one.

Module 4-FRESTM(Frequently restrictedly embedded subtree mining):

An Apriori-based data mining method, which progressively enumerates all candidate subtrees from a given set of unordered trees, level by level, using the rightmost expansion methods. In the initialization phase, frequent 1-subtrees and 2-subtrees are discovered first. To enumerate all frequent 1-subtrees, i.e., frequent single labels, we traverse every node of every tree to create an inverted index structure for each unique label appearing in the trees. Specifically, for each unique label, we maintain a list of IDs of supporting trees, denoted by STL, in which the label appears. By comparing its |STL| with the given minsup, we can decide whether the label is frequent or not.

Module 5-Generate Subtree:

The Generated Subtree to grow frequent subtrees level by level through pairwise joining and leg attachment methods. Notice that when $|FSTk|$ reaches zero, no more frequent $(k + 1)$ -subtrees can be generated and hence the discovering process terminates. Please notice that $|FSTk|$ can be as small one to allow self-joining and leg attachment operations.

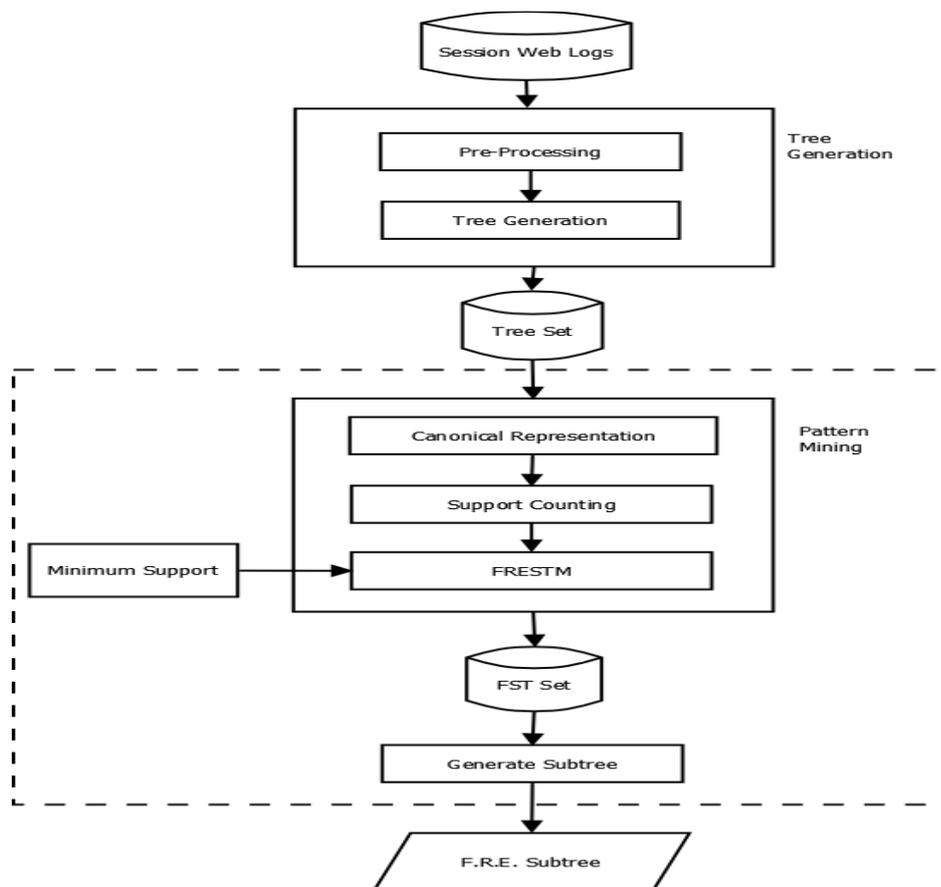


Figure :System Architecture

V. CONCLUSION

In this paper, we formalize a restrictedly embedded subtree mining problem, which has potential applications in many domains where data can be naturally represented as unordered unrooted trees. We study the properties of the canonical form of unordered unrooted trees and propose novel rightmost tree expansion techniques that can systematically, correctly, and efficiently generate all candidate subtrees. FRESTM, to solve the tree mining problem at hand. To the best of our knowledge, this is the first algorithm for finding restrictedly embedded subtree patterns in multiple unordered unrooted trees. Experimental results based on synthetic and real-world data demonstrate the good performance of the proposed algorithm.

REFERENCES

- [1] Sen Zhang, Zhihui Du, and Jason T. L. Wang, "New Techniques for Mining Frequent Patterns in Unordered Trees," *IEEE TRANSACTIONS ON CYBERNETICS*, VOL. 45, NO. 6, pp. 1113–1125, JUNE 2015.
- [2] M. H. Chehreghani, C. Lucas, and M. Rahgozar, "OInduced: An efficient algorithm for mining induced patterns from rooted ordered trees," *IEEE Trans. Syst., Man, Cybern. A, Syst. Humans*, vol. 41, no. 5, pp. 1013–1025, Sep. 2011.
- [3] S. Zhang and J. T. L. Wang, "Discovering frequent agreement subtrees from phylogenetic data," *IEEE Trans. Knowl. Data Eng.*, vol. 20, no. 1, pp. 68–82, Jan. 2008.
- [4] Y. Chi, Y. Xia, Y. Yang, and R. R. Muntz, "Mining closed and maximal frequent subtrees from databases of labeled rooted trees," *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 2, pp. 190–202, Feb. 2005.
- [5] M. J. Zaki, "Efficiently mining frequent trees in a forest: Algorithms and applications," *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 8, pp. 1021–1035, Aug. 2005.
- [6] K. G. Khoo and P. N. Suganthan, "Structural pattern recognition using genetic algorithms with specialized operators," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 33, no. 1, pp. 156–165, Feb. 2003.

- [7] C. H. Leung and C. Y. Suen, "Matching of complex patterns by energy minimization," *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 28, no. 5, pp. 712–720, Oct. 1998.
- [8] D. Shasha, J. T. L. Wang, K. Zhang, and F. Y. Shih, "Exact and approximate algorithms for unordered tree matching," *IEEE Trans. Syst., Man, Cybern.*, vol. 24, no. 4, pp. 668–678, Apr. 1994.
- [9] J. T. L. Wang, K. Zhang, K. Jeong, and D. Shasha, "A system for approximate tree matching," *IEEE Trans. Knowl. Data Eng.*, vol. 6, no. 4, pp. 559–571, Aug. 1994.
- [10] Y. Chi, Y. Yang, and R. R. Muntz, "HybridTreeMiner: An efficient algorithm for mining frequent rooted trees and free trees using canonical forms," in *Proc. 16th Int. Conf. Sci. Statist. Datab. Manage.*, Santorini Island, Greece, Jun. 2004.
- [11] A. Jiménez, F. Berzal, and J. Cubero, "POTMiner: Mining ordered, unordered, and partially-ordered trees," *Knowl. Inf. Syst.*, vol. 23, no. 2, May 2010, pp. 199–224.
- [12] M. L. Lee, L. H. Yang, W. Hsu, and X. Yang, "XClust: Clustering XML schemas for effective integration," in *Proc. 11th ACM Int. Conf. Inf. Knowl. Manage.*, McLean, VI, USA, Nov. 2002.
- [13] L. Liu and J. Liu, "Mining frequent embedded subtree from tree-like databases," in *Proc. Int. Conf. Internet Comput. Inf. Serv.*, Hong Kong, Sep. 2011.
- [14] S. Nijssen and J. N. Kok, "Efficient discovery of frequent unordered trees," in *Proc. 1st Int. Workshop Mining Graphs, Trees, Sequences(MGTS)*, 2003.
- [15] D. Shasha, J. T. L. Wang, and S. Zhang, "Unordered tree mining with applications to phylogeny," in *Proc. IEEE Int. Conf. Data Eng.*, Boston, MA, USA, pp. 708–719, 2004.
- [16] L. Zou, Y. Lu, H. Zhang, R. Hu, and C. Zhou, "Mining frequent induced subtree patterns with subtree-constraint," in *Proc. 6th IEEE Int. Conf. Data Mining (ICDM) Workshop*, Hong Kong, Dec. 2006.
- [17] T. Asai, H. Arimura, T. Uno, and S. Nakano, "Discovering frequent substructures in large unordered trees," in *Proc. 6th Int. Conf. Discov. Sci.*, Sapporo, Japan, Oct. 2003.