

# SS-ZG548: ADVANCED DATA MINING

# 14

# Mining with Tree Data



**Dr. Kamlesh Tiwari**  
Assistant Professor, Department of CSIS,  
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA

Oct 27, 2021 **ONLINE** (WILP @ BITS-Pilani July-Dec 2021)

<http://kti.warri.in/adm>

## Frequent Subtree Mining Problem

Given a tree database  $T_{db}$  and a minimum support threshold  $\sigma$ , find all subtrees that occur at least  $\sigma$  times in  $T_{db}$

### Support

- Transaction-Based Support: It is the number of transactions in  $T_{db}$  that contain at least one occurrence of subtree  $t$

$$\sum_{k_i \in T_{db}} d(t, k_i)$$

where  $d$  is an indicator variable for presence of subtree  $t$  in  $k_i$

- Occurrence-Based Support: Takes the repetition of items in a transaction and counts the subtree occurrences in the database as a whole

$$\sum_{k_i \in T_{db}} g(t, k_i)$$

where  $g$  denotes the total number of occurrences of subtree  $t$  in  $k_i$

- Hybrid-Match Support

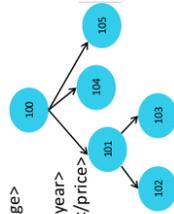
## Efficient Representation

'author(Giada De Laurentis)'

String	Index
book[category="COOKING"]	100
author	101
name[["Giada De Laurentis"]]	102
age["40"]	103
year["2005"]	104
price["30.00"]	105

```
<book category="COOKING">
<author>
```

```
<Name>Giada De Laurentis</Name>
<age>40</age>
</author>
<year>2005</year>
<price>30.00</price>
</book>
```



## Enumeration Methods

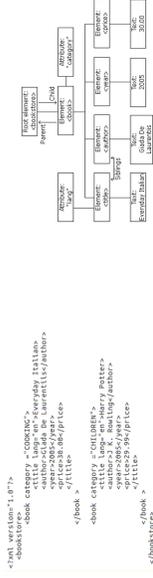
- Using join:  $ab + ac = abc$
- TreeMiner by Zaki 2005
- Tree Model Guided (TMG)

## Complex data forms<sup>1</sup>

Data can be categorized as 1) Transactional, 2) Relational, 3) Sequential, 4) Semi-structured, or 5) Unstructured

### XML: Extensible Markup Language

Data has a hierarchical relationship between nodes. NO loops



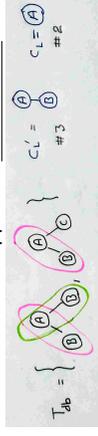
- Mine frequent patterns from XML (such as frequent trees)
- Can be modeled as a tree

<sup>1</sup>Book: Mining of Data with Complex Structures. By Fedja Hadzic, Henry Tan, Tharim S. Dillon

## AR Mining

Can we use apriory like algorithms (generate&test based approach)?

- NO for Occurrence-Based Support as anti-monotone property<sup>2</sup> violates



- YES for Transaction-Based Support

### Efficient Representation of XML tree

- Text format. No
- DFS string encoding: Traverse the tree in pre-order. Good for parent-child or ancestor-descendent or vertical relationships
- BFS string encoding: Traverse the tree in breadth-first order. Each sibling family is separated by the level coding

<sup>2</sup>Counter: In a tree database if  $\exists$  candidate subtrees  $C_L$  &  $C_R$ , where  $C_L \subseteq C_R$  s.t.  $C_L$  is frequent but  $C_R$  is infrequent, We say that  $C_L$  is pseudo frequent candidate subtree

## TreeMiner: Candidate Generation<sup>3</sup>

- Uses join approach: frequent k-subtrees are used to generate candidate (k+1)-subtrees
- Two k-subtrees X,Y are in the same prefix equivalence class iff they share a common prefix up to the (k-1)th node
- An element of a equivalence class with label x and attached to position  $i^{th}$  is denoted as  $(x, i)$
- An equivalence class with prefix P with elements  $(x, i), (y, j), (z, k)$  is denoted as  $[P] : (x, i), (y, j), (z, k)$
- Class Extension

- Let P be a prefix class with encoding P, and let  $(x, i)$  and  $(y, j)$  denote any two elements in the class
- Let  $[P_x]$  denote the class representing extensions of element  $(x, i)$

<sup>3</sup>TreeMiner: An Efficient Algorithm for Mining Embedded Ordered Frequent Trees. By Zaki, Mohammed J. AI Advanced Methods for Knowledge Discovery from Complex Data, pp 123-151, Springer, 2005

## TreeMiner: Candidate Generation

- Define a join operator  $\otimes$  on the two elements, denoted  $(x, i) \otimes (y, j)$  as follows

Case-1 ( $i = j$ ):

- If  $P \neq \phi$ , add  $(y, i)$  and  $(y, j + 1)$  to class  $[P_x]$
- If  $P = \phi$ , add  $(y, j + 1)$  to class  $[P_x]$

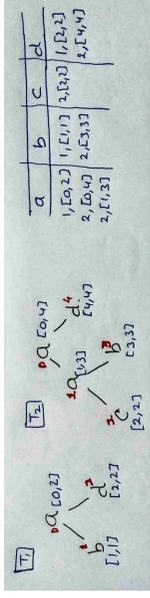
Case-2 ( $i > j$ ): add  $(y, i)$  to class  $[P_x]$

Case-3 ( $i < j$ ): no new candidate is possible in this case

- Enumerates all possible  $(k+1)$  subtrees with the common prefix P using  $k^{th}$  approaches, 2) Store only hyperlinks of subtrees in a tree database instead of creating a local copy for each generated subtree, 3) Subtree encoding is used to uniquely identify subtrees
- DFS is used to find scope of node  $n_i : [l, r]$  and scope-list that is a triple  $(t, m, s)$ . t: tree ID, m: match label, s: scope

Thank You!

## TreeMiner: Candidate Generation



- Two Type of Scope-list joins ( $s_x = [l_x, u_x]$ ,  $s_y = [l_y, u_y]$ )
  - In-scope Join: Used for ancestor-descendant or parent-child node extension.  $s_x \supset s_y$  iff  $l_x \leq l_y$  and  $u_y \leq u_x$
  - Out-scope Join: Used for sibling extension.  $s_x < s_y$  iff  $l_y \leq u_x$
- Add the triple  $(t_y, m_y \cup l_x, s_y)$

ab	acd	bcd
1,0,[1,1]	1,0,[2,2]	-
2,0,[2,3]	2,0,[4,4]	-
2,1,[3,3]		
	abd	
	1,0,[2,2]	
	2,0,[4,4]	

- Perform all possible joins. Frequency count of a subtree is equal to the size of its scope-list

Thank you very much for your attention!

Queries ?