

Efficient Combined Approach for Frequent Subgraph mining

By Md Mahamudur Rahaman Sajib

Outline

- Introduction
- Graph reminders
- Depth First Search (DFS) codes and tree
- Algorithm Specification

Introduction

- Extending APriori algorithms for itemsets and sequences to graphs
- Formulate new labeling method for easier graph testing that allows sorting of all graphs: DFS canonical label
- Use depth first search on hierarchical structure for faster performance, instead of breadth first search as in standard apriori algorithms

Graph basics

- It works on labeled simple graphs
- Labeled graph $G = (V, E, L, I)$
 - V set of vertices
 - $E \subseteq V \times V$ set of edges
 - L set of labels
 - $I : (V \cup E) \rightarrow L$ labeling of vertices and edges

Graph basics

- Definition: An isomorphism is a bijective function

$$f : V(G) \rightarrow V(H), (f(u), f(v)) \in E(H)$$

Goal

- Given dataset of graphs $GS = \{G_i \mid i=1..n\}$ and minimum support value, define

$$\xi(g, G) = \{ 1 \text{ if } g \text{ is isomorphic, } 0 \text{ otherwise } \}$$

$$\sigma(g, GS) = \sum \xi(g, G_i) \rightarrow \text{frequency of graph } g \text{ in } GS$$

- Frequent Subgraph Mining: – find graphs g in GS such that their frequency is greater or equal to minimum support

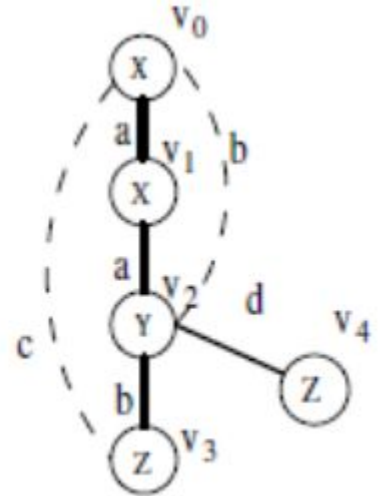
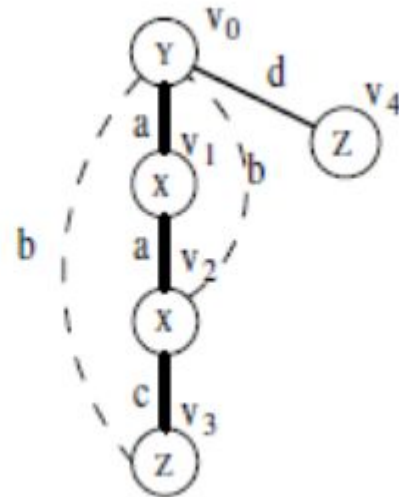
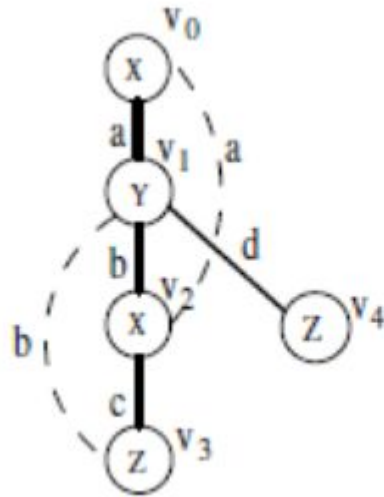
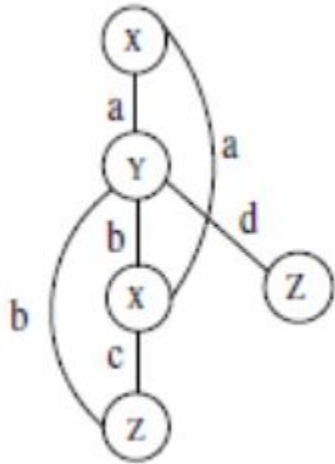
Idea Outline

- Instead of searching graphs and testing for isomorphism we construct canonical DFS codes
- Each graph has a canonical DFS code and the codes are equivalent if the graphs are isomorphic
- The codes are based on DFS trees

DFS tree

- Mark vertices in the the order they are traversed $v_i < v_j$ if v_i is traversed before v_j this constructs a DFS tree T , denoted GT
- DFS induces a linear order on vertices
- DFS divides edges in two sets – forward edge set: (v_i, v_j) where $v_i < v_j$ – backward edge set: (v_i, v_j) where $v_i > v_j$
- There are huge number of DFS trees for single graph

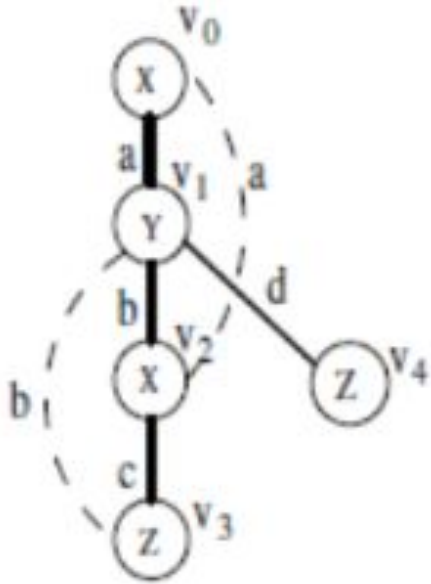
DFS tree



Linear orders

- A linear order of vertices defines a linear order of edges
 $(u, v) < (u, w)$ if $v < w$ and $(u, v) < (v, w)$ if $u < v$, $e_1 < e_2$
and $e_2 < e_3$ implies $e_1 < e_3$
- Linear order of edges is DFS code

Linear orders



$\{ (v_0, v_1), (v_1, v_2), (v_2, v_0), (v_2, v_3), (v_3, v_1), (v_1, v_4) \}$ is the linear order of edges for this dfs tree

DFS codes

- DFS code is a sequence of 4-tuples containing an edge and three labels
- Assume that there is an order on the labels
- This order together with the edge order defines an order for any two 4-tuples
- This extends to DFS code using a lexicographic encoding

DFS codes

- DFS code can be expanded with vertex and edge labels

Minimum DFS code

- Let the canonical DFS code to be the lexicographically smallest code that can be constructed from G (denoted $\min(G)$)
- Theorem: Given two graphs G and H , they are isomorphic if and only if $\min(G)=\min(H)$
- Subgraph mining: – Mining frequent subgraphs is equivalent mining their corresponding minimum DFS codes – Can be done sequentially by pattern mining algorithms

DFS Code Tree

- Definition: DFS code's parent and child

$$A = (a_0, a_1, a_2, \dots, a_M)$$

$$B = (a_0, a_1, a_2, \dots, a_M, b)$$

A is B's parent and B is child of A

- DFS Code Tree: – each node represents DFS code – relations between parents and children complies with previous definition – siblings are consistent with DFS lexicographic order

DFS Code Tree

- Theorem (frequency anti monotone): If a graph G is frequent, then any subgraph of G is frequent. If G is not frequent, then any graph which contains G is not frequent.

OR

- If a DFS code α is frequent, then every ancestor of α is frequent. If α is not frequent then every descendant of α is not frequent.

DFS Code Tree

- Some graphs can have more DFS nodes corresponding to it in DFS Code Tree
- The first occurrence is the minimum DFS code
- Theorem: If DFS code is not the minimum one, we can prune the entire subtree below this node, and still preserve DFS Code Tree Covering
- Pre-order searching of DFS Code Tree guarantees that we can enumerate all potential frequent subgraphs

Algorithm

GraphSet_projection(GS,FS)

- sort labels of the vertices and edges in GS by frequency;

- remove infrequent vertices and edges;

- relabel the remaining vertices and edges (descending);

- S1 := all frequent 1-edge graphs;

- sort S1 in DFS lexicographic order;

- FS := S1;

- for each edge e in S1

 - do init g with e

 - set Subgraph_mining(GS,FS,g);

 - GS := GS - e;

 - if |GS| < minSup break;

Algorithm

Subgraph_mining(GS,FS,g)

if $g \neq \min(g)$ return;

FS := FS \cup {g};

enumerate g in each graph in GS and count g's children;

for each c (child of g) do

if support(c) \geq minSup

Subgraph_mining(GS,FS,c);