# Subgraph Isomorphism

G1 is subgraph isomorphic to G2，denoted as G1 ⊆ G2:

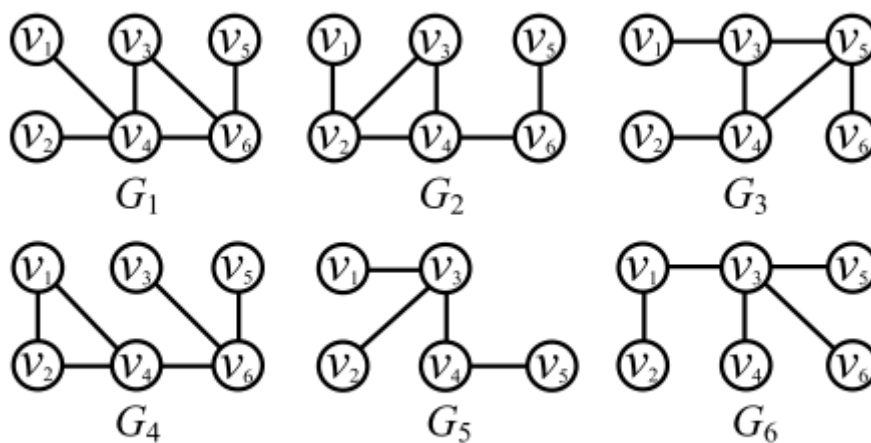iff there is an injective function f : V(G1) →V(G2), such that $\forall (u, v) \in E(G1)$, $(f(u),f(v)) \in E(G2)$

If G1 is subgraph isomorphic to G2,

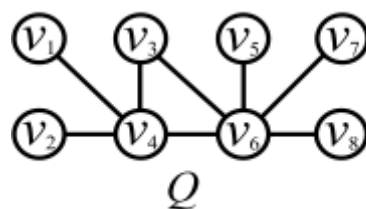G1 is called a *subgraph* of G2，G2 is called a *supergraph* of G1.

## Problem Statement

- a graph database D = {$G_1,G_2,..., G_n$}，each Gi ∈D is called a data graph
- a graph Q called *query-graph*
- *supergraph search* : find all graphs in D, that are subgraph isomorphic to Q
- Answer set *A(Q)*= {$G_i$|$G_i$ ∈D,$G_i$ ⊆ Q}

## Example



(a) Database Graphs

(b) Query Graph

Answer set A(Q)= {$G_1$,$G_5$,$G_6$}

## EXISTING SOLUTIONS

大多数现有的超图检索解决方案遵循修剪-验证框架，该框架根据修建阶段的特征修建错误答案（不符合特征的图数据），并在验证阶段对剩余图执行子图同构测试。

- 依靠频繁的子图挖掘算法来生成特征，开销昂贵，并且不能生成大的特征
- 验证阶段的开销也很昂贵
- 以固定的顺序处理特征而不考虑他们与查询图之间的关系

## contribution

- 提出了DGTree的索引结构。可以从图数据中直接选择特征图，而不用依靠频繁子图挖掘算法。并且在构建DGTree的过程中，考虑到了修剪和结构共享的能力。
- 基于特征树，提出了一种新的查询算法。设计了一个评分函数，评估特征图的共享成本和修剪能力，以便可以避免处理无用的特征。
- 提出了改进优化。

## DGTree: A Full-structure Index

## Match

a graph P with nodes {$u_1$,$u_2$,...,$u_{|V(P)|}$}

a data-graph G

a match  $f$  of P in G is a mapping from V(P) to V(G) such that the following two conditions hold:

- (**Conflict-free**) : for any pair of nodes $u_i \in$ V(P) and $u_j \in$ V(P) $(u_i \mathrel{!}= u_j)$, $f(u_i) \mathrel{!}= f(u_j)$
- (**Structure-preserved**): For any edge $(u_i, u_j) \in$ E(P), $(f(u_i), f(u_j)) \in$ E(G)

$f = [v_1, v_2, \ldots, v_{|V(P)|}]$ to denote the match f

$f(u_i) = v_i$ , for any $1 \leqslant i \leqslant |V(P)|$.

If $f(u_i) = v_i$, we have $f-1(v_i) = u_i$

## DGTree Structure

each tree-node  $g$ :

| $g$.children | The set of child tree-nodes of $g$. A leaf tree-node is a tree-node $g$ with $|g.\text{children}| = 0$. |
|---|---|
| $g$.graph | The feature-graph of $g$. The set of nodes in the feature-graph is represented by $\{1, 2, \ldots, |V(g.\text{graph})|\}$. |
| $g$.grow-edge | The edge added to $g$.graph from the feature-graph of the parent tree-node of $g$. |
| $g$.edge-type | The type of $g$.grow-edge, which is CLOSE if no new node is created in $g$.graph and OPEN if a new node is created in $g$.graph after adding $g$.grow-edge. |
| $g.\mathcal{S}$ | The set of data-graphs containing $g$.graph, i.e., $g.\mathcal{S} = \{G | G \in \mathcal{D}, g.\text{graph} \subseteq G\}$. If $g$ is a leaf tree-node, $g.\mathcal{S}$ contains the data-graph equaling to $g$.graph, and thus the union of $g.\mathcal{S}$ for all leaf tree-nodes $g$ is $\mathcal{D}$. |
| $g.\mathcal{M}(G_i)$ | The set of matches of $g$.graph in $G_i$ for each $G_i \in g.\mathcal{S}$. |
| $g.\mathcal{S}^*$ | $g.\mathcal{S}^* \subseteq g.\mathcal{S}$. If $g$ is the root tree-node, we have $g.\mathcal{S}^* = \mathcal{D}$. If $g$ is a leaf tree-node, we have $|g.\mathcal{S}^*| = 1$. For a non-leaf tree-node $g$, we have: (1) for any $g_i \in g.\text{children}$ and $g_j \in g.\text{children}$ $(g_i \neq g_j)$, $g_i.\mathcal{S}^* \cap g_j.\mathcal{S}^* = \emptyset$; and (2) $\bigcup_{g_i \in g.\text{children}} g_i.\mathcal{S}^* = g.\mathcal{S}^*$. In other words, $g_i.\mathcal{S}^*$ for all $g_i \in g.\text{children}$ form a disjoint cover of $g.\mathcal{S}^*$. |
| $g$.score | The score of $g$.graph, which is used to select the best edge to grow. |

**Algorithm 1**: DGTreeConstruct(database $\mathcal{D} = \{G_1, \ldots, G_n\}$)

1   $g_r \leftarrow$ a new tree-node;
2   $g_r.$graph $\leftarrow$ a single-edge graph;
3   $g_r.\mathcal{S} \leftarrow \mathcal{D}$; $g_r.\mathcal{S}^* \leftarrow \mathcal{D}$; $g_r.$grow-edge $= \emptyset$;
4   **for** $G_i \in \mathcal{D}$ **and** $(v, v') \in E(G_i)$ **do**
5     $\lfloor$   $g_r.\mathcal{M}(G_i) \leftarrow g_r.\mathcal{M}(G_i) \cup \{[v, v'], [v', v]\}$;

6   TreeGrow($g_r$);
7   **return** $g_r$;

8   **Procedure** TreeGrow(tree-node $g$)
9   $\mathcal{H} \leftarrow$ CandidateFeature($g$);
10   $\mathcal{C} \leftarrow g.\mathcal{S}^*$;
11   **while** $\mathcal{C} \neq \emptyset$ **do**
12     $g^+ \leftarrow$ BestFeature($\mathcal{H}, \mathcal{C}$);
13     **if** $|g^+.\mathcal{S}^*| > 1$ **then**
14       $g^+.$graph $\leftarrow$ a graph by adding $g.$grow-edge in $g.$graph;
15       TreeGrow($g^+$);

16     **else**   $g^+.$graph $\leftarrow$ the graph in $g.\mathcal{S}^*$; $g^+.\mathcal{S} \leftarrow g^+.\mathcal{S}^*$;
17     $g.$children $\leftarrow g.$children $\cup \{g^+\}$;
18     $\mathcal{C} \leftarrow \mathcal{C} \setminus g^+.\mathcal{S}^*$;

**Algorithm 2**: CandidateFeature(tree-node $g$)

1   $\mathcal{H} \leftarrow \emptyset$;
2   **for** data-graph $G \in g.\mathcal{S}^*$ **and** match $f \in g.\mathcal{M}(G)$ **do**
3     **for** $u_i \leftarrow 1$ **to** $|f|$ **and** $v \in \mathsf{Nbr}(f(u_i), G)$ **do**
4       **if** $v \in f$ **then** $u_j \leftarrow f^{-1}(v)$; $t \leftarrow$ CLOSE;
5       **else** $u_j \leftarrow |f| + 1$; $t \leftarrow$ OPEN;
6       **if** $u_j > u_i$ **and** $(u_i, u_j) \notin E(g)$ **then**
7         $g^+ \leftarrow \mathcal{H}.\mathsf{Find}((u_i, u_j))$;
8         **if** $g^+ = \emptyset$ **then**
9           $g^+ \leftarrow$ a new tree-node;
10          $g^+.\mathsf{grow\text{-}edge} \leftarrow (u_i, u_j)$;
11          $g^+.\mathcal{S}^* \leftarrow \{G\}$; $g^+.\mathsf{score} \leftarrow 0$;
12          $g^+.\mathsf{edge\text{-}type} \leftarrow t$;
13          $\mathcal{H}.\mathsf{Push}(g^+)$;
14         **else** $g^+.\mathcal{S}^* \leftarrow g^+.\mathcal{S}^* \cup \{G\}$;

15   **for** data-graph $G \in g.\mathcal{S}$ **and** match $f \in g.\mathcal{M}(G)$ **do**
16     **for** $u_i \leftarrow 1$ **to** $|f|$ **and** $v \in \mathsf{Nbr}(f(u_i), G)$ **do**
17       **if** $v \in f$ **then** $u_j \leftarrow f^{-1}(v)$;
18       **else** $u_j \leftarrow |f| + 1$;
19       **if** $u_j > u_i$ **and** $(u_i, u_j) \notin E(g)$ **then**
20         $g^+ \leftarrow \mathcal{H}.\mathsf{Find}((u_i, u_j))$;
21         **if** $g^+ \neq \emptyset$ **then**
22           $g^+.\mathcal{S} \leftarrow g^+.\mathcal{S} \cup \{G\}$;
23           **if** $g^+.\mathsf{edge\text{-}type} =$ OPEN **then**
24            $g^+.\mathcal{M}(G) \leftarrow g^+.\mathcal{M}(G) \cup \{[f, v]\}$;
25          **else** $g^+.\mathcal{M}(G) \leftarrow g^+.\mathcal{M}(G) \cup \{f\}$;

26   **for** $g^+ \in \mathcal{H}$ **do**
27     compute $g^+.\mathsf{score}$; $\mathcal{H}.\mathsf{Update}(g^+)$;
28   **return** $\mathcal{H}$;

---

**Algorithm 3**: BestFeature(heap $\mathcal{H}$, uncovered graphs $\mathcal{C}$)

1   $g^+ \leftarrow \mathcal{H}.\mathsf{Pop}()$;
2   **while** $g^+.\mathcal{S}^* \nsubseteq \mathcal{C}$ **do**
3     $g^+.\mathcal{S}^* \leftarrow g^+.\mathcal{S}^* \cap \mathcal{C}$;
4     **if** $g^+.\mathcal{S}^* \neq \emptyset$ **then** {compute $g^+.\mathsf{score}$; $\mathcal{H}.\mathsf{Push}(g^+)$};
5     $g^+ \leftarrow \mathcal{H}.\mathsf{Pop}()$;
6   **return** $g^+$;

## Score

- 为了在查询处理中分担计算成本，应使选择的g包含g.graph的数据图的数量最大化，也就是最大化|g.S|,但这个会生成冗余的特征图，因此选择最大化|g.S$^*$|

$$g.\,score_1 = |g.\,S^*|$$

- 如果一个数据图中的一个特征图平均匹配数很小，那么它不太可能包含在一个查询图中，因此这样一个特征图的剪枝能力很强

$$平均匹配数 = \sum_{G \in g.S} \frac{|g.\,M(G)|}{|g.\,S|}$$

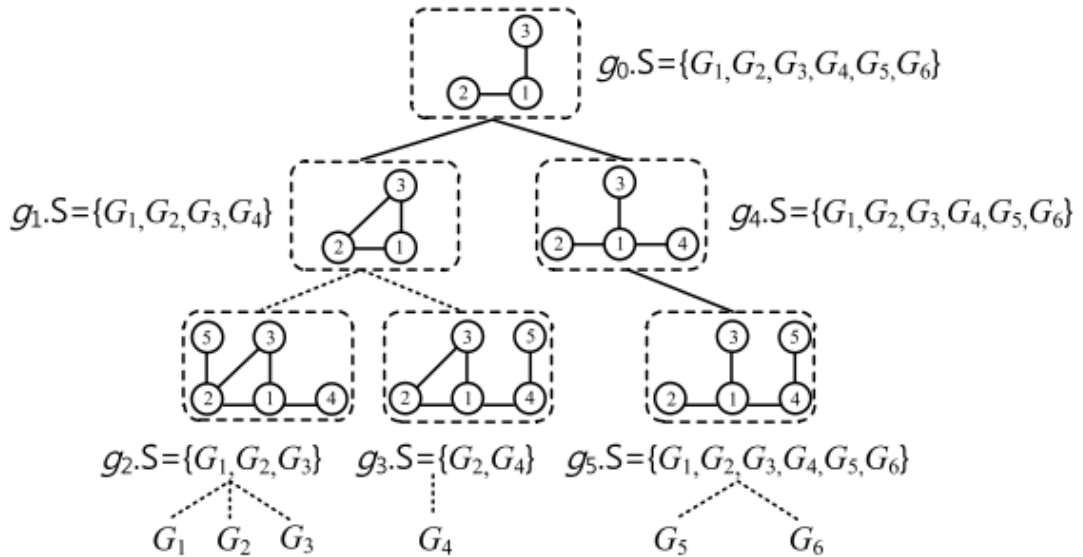$$g.\,score_2 = \frac{g.\,score_1}{平均匹配数} = \frac{|g.\,S^*| \times |g.\,S|}{\sum_{G \in g.S} |g.\,M(G)|}$$

## Example



Fig. 2: An Example of DGTree

# Query-dependent Supergraph Search

---

**Algorithm 4**: SuperGraphSearch(query $Q$, DGTree with root $g_r$)

---

1   $\mathcal{C} \leftarrow \{G | G \in g_r.\mathcal{S}, |E(G)| \leq |E(Q)|, |V(G)| \leq |V(Q)|\}$;

2   $\mathcal{H} \leftarrow \emptyset$; $\mathcal{A}(Q) \leftarrow \emptyset$;

3   $q \leftarrow$ a new entry;

4   $q.\text{tree-node} \leftarrow g_r$; $q.\mathcal{S}^* \leftarrow \mathcal{C}$; $q.\mathcal{M}(Q) \leftarrow \emptyset$;

5   **for** $(v, v') \in E(Q)$ **do**   $q.\mathcal{M}(Q) \leftarrow q.\mathcal{M}(Q) \cup \{[v, v'], [v', v]\}$;

6   compute $q.\text{score}$; $\mathcal{H}.\text{Push}(q)$;

7   **while** $\mathcal{C} \neq \emptyset$ **do**

8      $q \leftarrow \text{BestFeature}(\mathcal{H}, \mathcal{C})$; $g \leftarrow q.\text{tree-node}$;

9      **for** $g^+ \in g.\text{children}$ **do**

10         **if** $|g^+.\text{children}| = 0$ **then**

11            search a match $f$ of $g^+.\text{graph}$ by extending $q.\mathcal{M}(Q)$;

12            **if** $f \neq \emptyset$ **then**   $\mathcal{A}(Q) \leftarrow \mathcal{A}(Q) \cup g^+.\mathcal{S}$;

13            $\mathcal{C} \leftarrow \mathcal{C} \setminus g^+.\mathcal{S}$;

14         **else**   FeatureExpansion$(Q, q, g^+, \mathcal{H}, \mathcal{C})$;

15   **return** $\mathcal{A}(Q)$;

<br>

16   **Procedure** BestFeature(heap $\mathcal{H}$, candidate data-graph set $\mathcal{C}$)

17   $q \leftarrow \mathcal{H}.\text{Pop}()$;

18   **while** $q.\mathcal{S}^* \nsubseteq \mathcal{C}$ **do**

19      $q.\mathcal{S}^* \leftarrow q.\mathcal{S}^* \cap \mathcal{C}$;

20      **if** $q.\mathcal{S}^* \neq \emptyset$ **then**   $\{\text{compute } q.\text{score}; \mathcal{H}.\text{Push}(q)\}$;

21      $q \leftarrow \mathcal{H}.\text{Pop}()$;

22   **return** $q$;

---

**Algorithm 5**: FeatureExpansion(query-graph $Q$, entry $q$, tree-node $g^+$, heap $\mathcal{H}$, candidate data-graph set $\mathcal{C}$)

---

1   $q^+ \leftarrow$ a new entry;

2   $q^+.\text{tree-node} \leftarrow g^+$; $q^+.\mathcal{S}^* \leftarrow g^+.\mathcal{S} \cap \mathcal{C}$; $q^+.\mathcal{M}(Q) \leftarrow \emptyset$;

3   $(u_i, u_j) \leftarrow g^+.\text{grow-edge}$;

4   **for** match $f \in q.\mathcal{M}(Q)$ **do**

5      **if** $g^+.\text{edge-type} = \text{OPEN}$ **then**

6         **for** $v \in \text{Nbr}(f(u_i), Q)$ **do**

7            **if** $v \notin f$ **then**   $q^+.\mathcal{M}(Q) \leftarrow q^+.\mathcal{M}(Q) \cup \{[f, v]\}$;

8      **else if** $(f(u_i), f(u_j)) \in E(Q)$ **then**

9         $q^+.\mathcal{M}(Q) \leftarrow q^+.\mathcal{M}(Q) \cup \{f\}$;

10   **if** $q^+.\mathcal{M}(Q) \neq \emptyset$ **then**

11      compute $q^+.\text{score}$; $\mathcal{H}.\text{Push}(q^+)$;

12   **else**   $\mathcal{C} \leftarrow \mathcal{C} \setminus q^+.\mathcal{S}^*$;

## Query-dependent Feature Score

- 如果q的特征图包含在大量的候选数据图中，则可以最大限度地分担q的处理成本 ,所以最大化$|q.S^*|$

$$q.\,score_1 = |q.\,S^*|$$

- 在Q中具有少量匹配的特征图将具有较高的裁剪能力

$$q.\,score_2 = \frac{1}{q.\,M(Q)}$$

$$q.\,score = q.\,score_1 \times q.\,score_2$$