

# A Note on Restricted Forms of LGG

Ondřej Kuželka<sup>1</sup> and Jan Ramon<sup>2</sup>

<sup>1</sup> School of Computer Science & Informatics, Cardiff University, UK  
kuzelka0@cardiff.ac.uk

<sup>2</sup> Department of Computer Science, KU Leuven, Belgium  
jan.ramon@cs.kuleuven.be

**Abstract.** We study existence of a restricted least general generalization (LGG) with the property that LGGs of clauses from a pre-fixed set belong to this set. We show that there is no such LGG even in simple sets of clauses such as bounded-size clauses or treewidth-1 clauses.

## 1 Introduction

In this paper we study restricted forms of *least general generalization* (LGG) [6]. One such restricted form of LGG called *bounded LGG* was introduced in [4]. The main difference between ordinary LGG of some clauses  $A_1, A_2, \dots, A_k$  and their bounded LGG w.r.t. a set  $\mathcal{X}$  is that the latter type of LGG does not have to be the least general of all generalizations of these clauses, it merely suffices if it is less general than any other generalization from  $\mathcal{X}$ . In [3], it has been shown that bounded LGG can be used for hypothesis learning without having to resort to using exponential-time algorithms for  $\theta$ -subsumption if we allow the algorithm to potentially miss some hypotheses not from the set  $\mathcal{X}$  (e.g. some high-treewidth hypotheses). One property of bounded LGG was, however, still asking for a further study. When computing bounded LGG w.r.t. a set  $\mathcal{X}$ , it may often be the case that the resulting clause will not be from the set  $\mathcal{X}$ . For instance, when  $\mathcal{X}$  consists of clauses of treewidth bounded by  $k$ , it may be the case that a bounded LGG of some clauses w.r.t. this set will have treewidth higher than  $k$  even if the clauses to be generalized are all from the set  $\mathcal{X}$  as well. The question was whether there could be another type of restricted LGG which would not have this property, at least for some reasonable sets  $\mathcal{X}$ . We study this question in this paper and answer it negatively even for simple sets  $\mathcal{X}$ .

## 2 Preliminaries

A *directed graph* is a pair  $(V, E)$ , where  $V$  is a finite set of *vertices* and  $E \subseteq V \times V$  is a set of *edges*. Two vertices are said to be *adjacent* if they are connected by an edge. A graph  $H$  is *homomorphic* to a graph  $G$ , denoted  $H \preceq G$ , if there exists a mapping  $\varphi : V(H) \rightarrow V(G)$  such that if  $(u, v) \in E(H)$  then  $(\varphi(u), \varphi(v)) \in E(G)$ . Two graphs  $H$  and  $G$  are homomorphically equivalent if  $H \preceq G$  and  $G \preceq H$ . A graph  $G$  is a *core graph* if there is no strictly smaller graph homomorphically

equivalent to it. A *tree-decomposition* [7] of a graph  $G$ , denoted  $TD(G)$ , is a pair  $(T, \mathcal{X})$ , where  $T$  is a rooted unordered tree and  $\mathcal{X} = (X_z)_{z \in V(T)}$  is a family of subsets of  $V(G)$  satisfying: (i)  $\cup_{z \in V(T)} X_z = V(G)$ , (ii) for every  $\{u, v\} \in E(G)$ , there is a  $z \in V(T)$  such that  $u, v \in X_z$ , and (iii)  $X_{z_1} \cap X_{z_3} \subseteq X_{z_2}$  for every  $z_1, z_2, z_3 \in V(T)$  such that  $z_2$  is on the simple path connecting  $z_1$  with  $z_3$  in  $T$ . The set  $X_z$  associated with a node  $z$  of  $T$  is called the *bag* of  $z$ . The treewidth of  $TD(G)$  is  $\max_{z \in V(T)} |X_z| - 1$ , and the *treewidth* of  $G$ , denoted  $\text{tw}(G)$ , is the minimum treewidth over all tree-decompositions of  $G$ . By graphs of bounded treewidth we mean graphs of treewidth at most  $k$ , where  $k$  is some constant. For example, all trees have treewidth 1, cycles have treewidth 2, rectangular  $n \times n$  grids have treewidth  $n$ . A graph with treewidth 1 is a forest, possibly with loops.

A first-order-logic *clause* is a universally quantified disjunction of first-order-logic literals. For convenience, we do not write the universal quantifiers explicitly. We treat clauses as disjunctions of literals and as sets of literals interchangeably. To denote the number of literals in a clause  $A$ , we use the set notation  $|A|$ . A clause  $A$   $\theta$ -subsumes a clause  $B$  (denoted by  $A \preceq_{\theta} B$ ), if and only if there is a substitution  $\theta$  such that  $A\theta \subseteq B$ . If  $A \preceq_{\theta} B$  and  $B \preceq_{\theta} A$ , we call  $A$  and  $B$   $\theta$ -equivalent (written  $A \approx_{\theta} B$ ). A clause  $C$  is  $\theta$ -reducible if there exists a clause  $C'$  such that  $C' \approx_{\theta} C$  and  $|C'| < |C|$ . A clause with minimal number of literals  $\theta$ -equivalent to a clause  $C$  is its  $\theta$ -reduction.  $\theta$ -subsumption corresponds to homomorphism and  $\theta$ -reduction corresponds to core of a graph. The *Gaifman* graph of a clause  $A$  is the graph with one vertex for each variable  $v \in \text{vars}(A)$  and an edge for every pair of variables  $u, v \in \text{vars}(A)$ ,  $u \neq v$  such that  $u$  and  $v$  appear in a literal  $l \in A$ . The *treewidth of a clause* is equal to the treewidth of its Gaifman graph.  $\theta$ -subsumption and  $\theta$ -reduction can be computed in polynomial time for clauses which have bounded-treewidth  $\theta$ -reductions.

A clause  $C$  is said to be a least general generalization of clauses  $A$  and  $B$  (denoted by  $C \in \text{LGG}(A, B)$ ) if and only if  $C \preceq_{\theta} A$ ,  $C \preceq_{\theta} B$  and for every clause  $D$  such that  $D \preceq_{\theta} A$  and  $D \preceq_{\theta} B$  it holds  $D \preceq_{\theta} C$ . An LGG of two clauses  $C$ ,  $D$  can be computed in time  $\mathcal{O}(|C| \cdot |D|)$ . LGG can be used as an operator in the process of searching for hypotheses [1,5]. A problem of approaches based on least general generalization is that the size of an LGG of a set of examples can grow exponentially in the number of examples. In order to keep the LGGs reasonably small,  $\theta$ -reduction is typically applied on the result of each LGG iteration [1]. An alternative to LGG capable of exploiting tractability of restricted hypothesis classes, called *bounded LGG*, was introduced in [4]. We discuss bounded LGG in the next section w.r.t. its relationship to *LGG in a set*.

### 3 Bounded LGG and LGG in a Set

The concept of bounded LGG was introduced in [4] in order to exploit existence of hypothesis classes with tractable  $\theta$ -subsumption for learning based on LGG.

**Definition 1 (Bounded LGG).** *Let  $\mathcal{X}$  be a set of clauses. A clause  $B$  is said to be a bounded LGG of clauses  $A_1, A_2, \dots, A_n$  w.r.t. the set  $\mathcal{X}$  (denoted by*

$B \in \text{LGG}_{\mathcal{X}}(A_1, A_2, \dots, A_n)$  if and only if  $B \preceq_{\theta} A_i$  for all  $i \in \{1, \dots, n\}$  and if for every other clause  $C \in \mathcal{X}$  such that  $C \preceq_{\theta} A_i$  for all  $i \in \{1, \dots, n\}$ , it holds  $C \preceq_{\theta} B$ .

Note that neither the clauses to be generalized, nor the resulting bounded LGG w.r.t.  $\mathcal{X}$  are required to belong to the set  $\mathcal{X}$ . In fact, there are cases where we can show easily that there is no bounded LGG belonging to the set  $\mathcal{X}$ . The following example from [4] is one such case.

*Example 1.* Let  $\mathcal{X} = \{C_1, C_2, \dots\}$  be a set of clauses of the following form:  $C_1 = e(A_1, A_2)$ ,  $C_2 = e(A_1, A_2) \vee e(A_2, A_3)$ ,  $C_3 = e(A_1, A_2) \vee e(A_2, A_3) \vee e(A_3, A_4)$ , etc. Let us also have the following two clauses:  $A = e(X, Y) \vee e(Y, X)$  and  $B = e(X, Y) \vee e(Y, Z) \vee e(Z, X)$ . We would like to find a clause from  $\mathcal{X}$  which would be their LGG but this is impossible for the following reason. Any clause from  $\mathcal{X}$   $\theta$ -subsumes both  $A$  and  $B$  but none of them is least general because for any  $C_i \in \mathcal{X}$  we have  $C_{i+1} \not\preceq_{\theta} C_i$ ,  $C_{i+1} \preceq_{\theta} A$  and  $C_{i+1} \preceq_{\theta} B$ . On the other hand, bounded LGG, as actually defined, always exists which follows trivially from the fact that the conventional LGG as computed by Plotkin's algorithm [6] is also a bounded LGG. Nevertheless, it does not belong to the set  $\mathcal{X}$ .

Notice that the clauses  $A$  and  $B$  in the above example do not belong to the set  $\mathcal{X}$ . In fact, one can verify easily that for all the clauses from the set  $\mathcal{X}$  from the above example, there is always an LGG belonging to the set  $\mathcal{X}$ . Thus, one might conjecture that, in general, if we restrict the clauses of interest, which we may want to generalize, to be from the set  $\mathcal{X}$  then we will always be able to find a bounded LGG from the set  $\mathcal{X}$ <sup>3</sup>. This motivates the definition of the following, arguably quite natural, type of LGG which is studied in this paper.

**Definition 2 (LGG in a set  $\mathcal{X}$ ).** Let  $\mathcal{X}$  be a set of clauses. A clause  $B \in \mathcal{X}$  is said to be an LGG of clauses  $A_1, A_2, \dots, A_n \in \mathcal{X}$  in the set  $\mathcal{X}$  (denoted by  $B \in \text{LGG}_{\mathcal{X}}^{\text{in}}(A_1, A_2, \dots, A_n)$ ) if and only if  $B \preceq_{\theta} A_i$  for all  $i \in \{1, \dots, n\}$  and if for every other clause  $C \in \mathcal{X}$  such that  $C \preceq_{\theta} A_i$  for all  $i \in \{1, \dots, n\}$ , it holds  $C \preceq_{\theta} B$ .

There are several important differences between LGG in a set  $\mathcal{X}$  ( $\text{LGG}_{\mathcal{X}}^{\text{in}}$ ) and bounded LGG w.r.t. a set  $\mathcal{X}$  ( $\text{LGG}_{\mathcal{X}}$ ). Most importantly, bounded LGG w.r.t. a set  $\mathcal{X}$  is not required to belong to the set  $\mathcal{X}$  which is the property guaranteeing that it always exists. Since LGG in a set  $\mathcal{X}$  must belong to  $\mathcal{X}$ , it may be the case that it does not exist. Arguably for the sets  $\mathcal{X}$  in which  $\text{LGG}_{\mathcal{X}}^{\text{in}}$  exists, it would be preferable over  $\text{LGG}_{\mathcal{X}}$ , especially for the sets  $\mathcal{X}$  for which tractable  $\theta$ -subsumption algorithms exist (e.g. bounded-size or bounded-treewidth clauses). That is one of the reasons why, in this paper, we are interested in the question of existence of  $\text{LGG}_{\mathcal{X}}^{\text{in}}$  in several such sets of clauses. Here, we note that if  $A_1, A_2, \dots, A_n \in \mathcal{X}$  and  $\text{LGG}(A_1, A_2, \dots, A_n) \notin \mathcal{X}$  then this does not yet mean that  $\text{LGG}_{\mathcal{X}}^{\text{in}}(A_1, A_2, \dots, A_n)$  does not exist<sup>4</sup>.

<sup>3</sup> Actually, the main negative results presented in this paper show that this is not the case in the majority of interesting cases.

<sup>4</sup> If this was the case then the problem of existence of  $\text{LGG}_{\mathcal{X}}^{\text{in}}$  would be almost trivial.

The results of Horváth and Turán [2] imply that an  $\text{LGG}_{\mathcal{X}}^{\text{in}}$  operator exists in the class of forests of rooted directed trees (although not using this terminology). The results presented in this paper actually show that, a bit surprisingly, the results of Horváth and Turán cannot be extended much (for instance they cannot be generalized to the class of treewidth-1 graphs).

## 4 No LGGs in Sets of Bounded-Size Clauses

In this section, to start with a simpler problem before we tackle the question of existence of an LGG operator in the set of bounded-treewidth clauses, we consider the question of existence of LGG in the set of clauses consisting of at most  $k$  atoms. We show that there is no  $\text{LGG}_{\mathcal{X}}^{\text{in}}$  operator in the sets of clauses consisting of at most  $k$  atoms where  $k$  is an integer greater or equal to 4.

**Theorem 1.** *If  $n \geq 4$  then there is no LGG operator in the set  $\mathcal{X}_n$  of clauses with at most  $n$  atoms based on one binary predicate. There is an LGG operator in the set  $\mathcal{X}_3$  of clauses consisting of at most 3 binary atoms with the same predicate.*

The next example shows<sup>5</sup> that there are clauses  $A, B \in \mathcal{X}_3$  such that  $\text{LGG}(A, B) \cap \mathcal{X}_3 = \emptyset$  and  $\text{LGG}_{\mathcal{X}_3}^{\text{in}}(A, B) \neq \emptyset$ .

*Example 2.* Let us have the following two clauses:  $A = e(X, Y) \vee e(Y, X)$  and  $B = e(X, Y) \vee e(Y, Z) \vee e(Z, X)$ . Their *conventional*  $\theta$ -reduced LGG is  $\text{LGG}(A, B) = e(X_1, X_2) \vee e(X_2, X_3) \vee e(X_3, X_4) \vee e(X_4, X_5) \vee e(X_5, X_6) \vee e(X_6, X_1)$ , and thus  $\text{LGG}(A, B) \cap \mathcal{X}_3 = \emptyset$ . However, there exists an  $\text{LGG}_{\mathcal{X}_3}^{\text{in}}(A, B)$ , for instance,  $e(W, X) \vee e(X, Y) \vee e(Y, Z) \in \text{LGG}_{\mathcal{X}_3}^{\text{in}}(A, B)$ .

Along the same lines, we can show that if we allow more than one binary predicate, the situation becomes even worse.

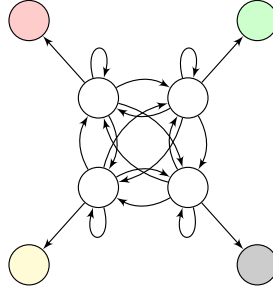
**Theorem 2.** *If  $n \geq 3$  then there is no LGG operator in the set  $\mathcal{X}_n^{(2)}$  of clauses with at most  $n$  atoms with two different binary predicates.*

We could see in this section, which was mostly meant to illustrate the general problem of existence of LGGs in sets of clauses, that size of the clauses is not a very good measure for defining sets of clauses with an LGG operator. This is a bit unfortunate but not very surprising.

## 5 No LGGs in Sets of Treewidth-1 Clauses

In this section, we show that there is no LGG in the set of clauses with treewidth 1. This is quite surprising given the positive result of Horváth and Turán [2]. However, there is no disagreement between this positive result and our negative result

<sup>5</sup> From this, it also follows that in order to show that, in general, there is no LGG in the set  $\mathcal{X}_3$ , it is not enough to show that  $\theta$ -reduced Plotkin's LGG of some two clauses from  $\mathcal{X}_4$  has more than 4 atoms.



**Fig. 1.** The graph corresponding to clause  $C \in LGG(A, B)$  from the proof of Theorem 3.

as the negative result depends on the fact that graphs with loops have treewidth 1 too whereas loops are not allowed in the other setting corresponding to the positive result.

**Theorem 3.** *There is no LGG operator for the set of clauses with treewidth 1.*

*Proof.* Let us have two clauses

$$\begin{aligned}
 A = & \text{red}(a1) \vee \text{green}(a2) \vee \text{yellow}(a3) \vee \text{black}(a4) \vee e(a5, a1) \vee a(a5, a2) \vee \\
 & \vee e(a5, a6) \vee e(a6, a5) \vee e(a6, a3) \vee e(a6, a4) \vee e(a5, a5) \vee e(a6, a6) \\
 B = & \text{red}(b1) \vee \text{yellow}(b2) \vee \text{green}(b3) \vee \text{black}(b4) \vee e(b5, b1) \vee a(b5, b2) \vee \\
 & \vee e(b5, b6) \vee e(b6, b5) \vee e(b6, b3) \vee e(b6, b4) \vee e(b5, b5) \vee e(b6, b6)
 \end{aligned}$$

The conventional LGG of the clauses  $A$  and  $B$  is a clause  $C$  which represents the graph shown in Figure 1. The clause  $C$  is not  $\theta$ -reducible (i.e. the corresponding labeled graph is a core) and has treewidth greater than 1 as it contains a clique on 4 vertices. As we have already explained this does not guarantee that there is no LGG of  $A$  and  $B$  in the set of clauses of treewidth 1. We therefore need to prove that there is indeed no such clause of treewidth 1, which we will do by contradiction.

Let us assume that there is a clause  $D$  which is an LGG of  $A$  and  $B$  and which has treewidth 1. Such a clause must correspond to a labeled tree or forest, possibly with loops. It follows from the definitions of LGG and LGG in a set that  $D$  must also  $\theta$ -subsume  $C$ . Let us define a family of clauses

$$\begin{aligned}
 E_1 = & \text{red}(Y_1) \vee e(X_1, Y_1) \vee e(X_1, X_2) \vee \text{green}(Y_2) \vee e(X_2, Y_2) \vee e(X_2, X_3) \vee \\
 & \vee \text{black}(Y_3) \vee e(X_3, Y_3) \vee e(X_3, X_4) \vee \text{yellow}(Y_4) \vee e(X_4, Y_4)
 \end{aligned}$$

$$\begin{aligned}
E_2 = & \text{red}(Y_1) \vee e(X_1, Y_1) \vee e(X_1, X_2) \vee \text{green}(Y_2) \vee e(X_2, Y_2) \vee e(X_2, X_3) \vee \\
& \vee \text{black}(Y_3) \vee e(X_3, Y_3) \vee e(X_3, X_4) \vee \text{yellow}(Y_4) \vee e(X_4, Y_4) \vee \\
& \vee e(X_4, X_5) \vee \text{red}(Y_5) \vee e(X_5, Y_5) \vee \cdots \vee \text{yellow}(X_8) \vee e(X_8, Y_8) \\
& \cdots \\
E_k = & \text{red}(Y_1) \vee e(X_1, Y_1) \vee e(X_1, X_2) \vee \cdots \vee \text{yellow}(Y_{4k}) \vee e(X_{4k}, Y_{4k}). \\
& \cdots
\end{aligned}$$

Clearly, each  $E_i$   $\theta$ -subsumes  $C$ . By the assumption that  $D$  is an LGG in the set of clauses of treewidth 1, each  $E_i$  should also  $\theta$ -subsume  $D$  (because each  $E_i$  has treewidth 1 and  $\theta$ -subsumes  $A$  and  $B$ ). Let us denote  $D_i = E_i\theta$  where  $E_i\theta_i \subseteq D$  and  $\theta_i$  is an arbitrary suitable substitution. Since  $D$   $\theta$ -subsumes the clause  $C$ , no vertex in the graph corresponding to the clause  $D$  can be adjacent to two vertices labeled by different colors, i.e. the clause  $D$  cannot contain simultaneously e.g. literals  $e(X, Y)$ ,  $e(X, Z)$ ,  $\text{yellow}(Y)$  and  $\text{red}(Z)$ . It follows that if  $e(X_j, X_{j+1}) \vee e(X_{j+1}, X_{j+2}) \subseteq E_i$  then  $\theta_i$  cannot map  $X_j$  and  $X_{j+2}$  on the same term (this follows from the construction of  $E_i$ 's). For similar reasons,  $\theta_i$  cannot map  $X_j$  and  $X_{j+1}$  to the same term. Since  $D$  corresponds to a directed tree, possibly with loops or cycles of length 2, and therefore contains no simple cycles of length greater than 2, it follows that no two  $X_j \neq X_{j'}$  can be mapped to the same term in  $D$ . However, since  $D$  is finite, there must be  $E_k$  such that  $E_k \preceq_{\theta} A$ ,  $E_k \preceq_{\theta} B$  but  $E_k \not\preceq_{\theta} D$  which is a contradiction with  $D$  being an LGG of  $A$  and  $B$  in the set of treewidth-1 clauses. It follows that there is no finite LGG of  $A$  and  $B$  in this set of clauses.  $\square$

Note that the above theorem shows the existence of a counterexample only for treewidth-1 clauses and not for treewidth- $k$  clauses in general. Thus, theoretically, it might be the case that there is an LGG operator in the class of clauses of treewidth at most  $k$ , where  $k > 1$ , and a proof would still be needed to disprove such a conjecture for general  $k$ . This seems unlikely, though.

## 6 Conclusions

The problems studied in this paper were motivated by the question whether bounded LGG w.r.t. a set  $\mathcal{X}$ , introduced in [4], could not be replaced by another type of LGG guaranteeing that the resulting generalized clauses would belong to the set  $\mathcal{X}$ , at least when generalizing clauses from  $\mathcal{X}$ . We have shown that such an alternative LGG does not exist already for natural and simple sets  $\mathcal{X}$  such as the set of bounded-size clauses and the set of treewidth-1 clauses. Thus, to our best knowledge, bounded LGG remains the only candidate for an LGG capable of exploiting tractability of bounded-treewidth clauses for learning based on LGG.

**Acknowledgement.** This work was supported by ERC Starting Grant 240186 “MiGraNT: Mining Graphs and Networks, a Theory-based approach”. The first author is supported by a grant from the Leverhulme Trust (RPG-2014-164).

## References

1. T. Horváth, G. Paass, F. Reichartz, and S. Wrobel. A logic-based approach to relation extraction from texts. In *ILP*, pages 34–48, 2009.
2. T. Horváth and G. Turán. Learning logic programs with structured background knowledge. *Artif. Intell.*, 128(1-2):31–97, 2001.
3. O. Kuželka. *Fast Construction of Relational Features for Machine Learning*. PhD thesis, CTU in Prague, 2013.
4. O. Kuželka, A. Szabóová, and F. Železný. Bounded least general generalization. In *ILP 2012*, pages 116–129, 2012.
5. S. Muggleton and C. Feng. Efficient induction of logic programs. In *ALT*, pages 368–381, 1990.
6. G. Plotkin. *A note on inductive generalization*. Edinburgh University Press, 1970.
7. N. Robertson and P. D. Seymour. Graph minors .xiii. the disjoint paths problem. *J. Comb. Theory, Ser. B*, 63(1):65–110, 1995.