

A similarity measure for graphs with low computational complexity

Matthias Dehmer ^{a,*}, Frank Emmert-Streib ^b, Jürgen Kilian ^a

^a Technische Universität Darmstadt, 64289 Darmstadt, Germany

^b Stowers Institute for Medical Research, 1000 E. 50th Street, Kansas City, MO 64110, USA

Abstract

We present and analyze an algorithm to measure the structural similarity of generalized trees, a new graph class which includes rooted trees. For this, we represent structural properties of graphs as strings and define the similarity of two graphs as optimal alignments of the corresponding property strings. We prove that the obtained graph similarity measures are so called Backward similarity measures. From this we find that the time complexity of our algorithm is polynomial and, hence, significantly better than the time complexity of classical graph similarity methods based on isomorphic relations. © 2006 Elsevier Inc. All rights reserved.

Keywords: Graph theory; Graph similarity; Dynamic programming; Computational complexity

1. Introduction

In this paper, we examine the computational complexity of a novel algorithm to measure the structural similarity of graphs. Thereby, we deal with a special graph class: directed and hierarchical graphs which we call *generalized trees*. There are classical approaches to measure the structural similarity of graphs, e.g., [4,11–13,20,25]. The typical and classical approach to measure the structural similarity of graphs is based on isomorphic relations [11–13,20,25]. For example, Zelinka [25] was the first who measured the distance between isomorphism classes of unlabeled graphs. The graph similarity measure of Zelinka is based on the principle that two graphs are more similar the bigger the common induced isomorphic subgraph is. Furthermore Sobik [20] and Kaden [12] generalized the measure of Zelinka for arbitrary and labeled graphs of different orders.¹ The main disadvantage of these graph measures is that the computational complexity is \mathcal{NP} -complete [20,23]. Our approach to measure the structural similarity of generalized trees is completely different from similarity or distance measures which are based on isomorphic relations mentioned above. Our novel similarity measure is

* Corresponding author.

E-mail addresses: dehmer@informatik.tu-darmstadt.de (M. Dehmer), fes@stowers-institute.org (F. Emmert-Streib), kilian@noteserver.org (J. Kilian).

¹ If we denote the given graph as $G = (V, E)$, then the order of these graphs is defined by $|V|$.

mainly based on a new representation of graphs as strings, whose components represent structural properties of the graph. We call these strings *property strings* of a generalized tree. In Section 3, the structural similarity of two graphs will be defined as the optimal alignment of the underlying property strings. In the present paper we examine the computational complexity of our algorithm. We want to emphasize, that generalized trees are more complex than normal directed rooted trees. Therefore this graph class can represent more complex problems than ordinary trees as demonstrated, e.g., in [26]. In [6], Emmert-Streib et al. showed that our novel method is applicable to measure the structural similarity of a more general graph type (unweighted and undirected graphs) by a local decomposition of the graph in generalized trees. Due to [6], considering generalized trees is not necessarily a restriction in terms of generality.

This paper is organized as follows: In Section 2 we give a review of classical graph and tree similarity methods to distinguish these methods better from our novel method. We repeat in Section 3 the construction of our graph similarity measures and prove for the first time that there are *Backward similarity measures*. The examination of the computational complexity of our method will be presented in Section 4. Section 5 finishes this paper with a summary.

2. Measuring the structural similarity of graphs

2.1. Preliminaries

In this section we give first some fundamental definitions of graph theory and then review some approaches to measure the structural similarity of trees and graphs. In order to distinguish trees from generalized trees we need the following definition.

Definition 2.1. Let $\mathcal{H} = (\hat{V}, E_1)$ be an directed rooted tree. The vertex set \hat{V} is defined by

$$\hat{V} := \{v_{0,1}, v_{1,1}, v_{1,2}, \dots, v_{1,\sigma_1}, v_{2,1}, v_{2,2}, \dots, v_{2,\sigma_2}, \dots, v_{h,1}, v_{h,2}, \dots, v_{h,\sigma_h}\},$$

where $v_{i,j}$ denotes the j th vertex on the i th level, $0 \leq i \leq h$, $1 \leq j \leq \sigma_i$. h denotes the depth of \mathcal{H} and σ_i is the number of vertices on level i . The edge set $\hat{E} := \hat{E}_1 \cup \hat{E}_2 \cup \hat{E}_3 \cup \hat{E}_4$ is defined as [16]

- (\hat{E}_1) forms the edge set of the underlying directed rooted tree \mathcal{H} .
- (\hat{E}_2) *Up-edges* associate analogously vertices of the tree hierarchy with one of their (dominating) predecessor vertices.
- (\hat{E}_3) *Down-edges* associate vertices of the tree hierarchy with one of their (dominated) successor vertices in terms of that tree hierarchy.
- (\hat{E}_4) *Cross-edges* associate vertices of the tree hierarchy, none of which is an (immediate) predecessor of the other in terms of the tree hierarchy.

If \hat{E}_2 , \hat{E}_3 or $\hat{E}_4 \neq \emptyset$ then $\hat{\mathcal{H}} = (\hat{V}, \hat{E})$ is called a generalized tree (Fig. 1).

Definition 2.2. Let $\mathcal{H} := (V, E), E \subseteq V \times V$ be an unlabeled, directed graph. $\mathcal{G} := (\hat{V}, \hat{E})$ with $\hat{V} \subseteq V$ and $\hat{E} \subseteq E$ denotes the partial graph of \mathcal{H} , $\mathcal{G} \subseteq \mathcal{H}$.

Definition 2.3. Let $\mathcal{H} := (V, E)$ be a directed graph and $\mathcal{G} := (\hat{V}, \hat{E})$ be a partial graph of \mathcal{H} . Moreover, it holds $\hat{E} = E \cap (\hat{V} \times \hat{V})$. Then we call \mathcal{G} the induced subgraph of \mathcal{H} .

In the following, we give the definition of graph isomorphism that provides the structural equivalence of two graphs [21].

Definition 2.4. Let $\mathcal{H} := (V, E)$ and $\mathcal{G} := (\hat{V}, \hat{E})$ be directed graphs. We call \mathcal{H} and \mathcal{G} isomorphic ($\mathcal{H} \cong \mathcal{G}$) if it exists a bijective mapping $\phi : V \rightarrow \hat{V}$ such that

$$E \ni (v_i, v_j) : \iff (\phi(v_i), \phi(v_j)) \in \hat{E}.$$

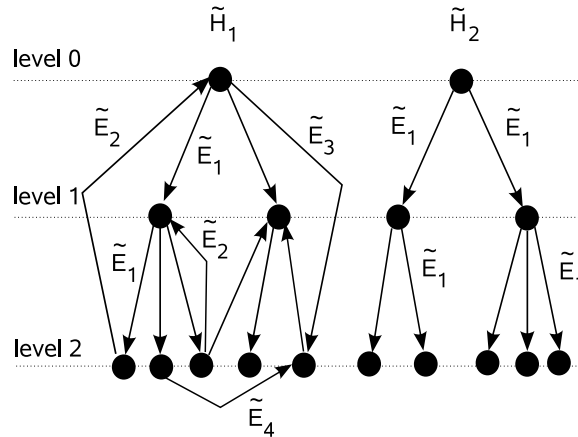


Fig. 1. $\tilde{\mathcal{H}}_1$ shows a generalized tree and his edge types fulfills Definition 2.1. In contrast $\tilde{\mathcal{H}}_2$ represents an ordinary directed rooted tree, which consists only of edges $e \in \tilde{E}_1$. An edge $e \in \tilde{E}_1$ over jumps always just one level, $e \in \tilde{E}_3$ over jumps at least one level, $e \in \tilde{E}_4$ does not necessarily over jump a level.

The mapping ϕ is called the isomorphism of \mathcal{H} on \mathcal{G} . Explaining informally two graphs are isomorphic iff, we obtain the first graph by renaming vertex labels of the second graph.

For completing this Section we state the definition of directed rooted trees.

Definition 2.5. We call $T = (V_T, E_T)$ a directed rooted tree iff E_T consists only of edges from E_1 (see Definition 2.1).

2.2. Classical tree similarity measures

There are known approaches [17,18] to measure the *edit distance* between strings $s_1, s_2 \in A^\star$ on the basis of string alignments, where A^\star is an arbitrary alphabet. Since one can consider special types of graphs as a generalization of the string concept, Tai [22] and Selkow [19] extended the string alignment problem to undirected and *ordered trees*. Thereby, we call a tree *ordered tree*, if the sequence of child vertices $v_1, v_2, \dots, v_{\delta(v)}$ of $v \in V_T$ is significant, where $\delta(v)$ denotes the degree of v . Höchstmann et al. [9] provide a method for measuring the structural similarity of RNA *secondary structures* representing a *forest* [9] of ordered trees. In [9] a forest is defined by a finite sequence of ordered trees. The approach of Höchstmann et al. is based on the tree alignment technique of Jiang et al. [10]. Further developments and novel applications of tree similarity measures can be found in [7,15,24].

2.3. Computational complexity of tree similarity measures

Assuming the vertex labeled, ordered trees $T_1 = (V_1, E_1), T_2 = (V_2, E_2)$, the basic *tree edit problem* [19,22] can be solved in time $\mathcal{O}(|V_1| \cdot |V_2| \cdot D_1^2 \cdot D_2^2)$, where D_1, D_2 are the associated maximum depths of T_1, T_2 . An alternative to tree edit has been stated by Jiang et al. [10]. In [10], Jiang et al. introduced an algorithm to measure the structural similarity of ordered and undirected trees with time complexity $\mathcal{O}(|V_1| \cdot |V_2| \cdot \deg(T_1) \cdot \deg(T_2))$, where $\deg(T_k) := \max(\{\delta(v_i^k) | v_i^k \in V_k, k = 1, 2\})$. The algorithm of Höchstmann et al. [9] has the time complexity

$$\mathcal{O}(|F_1| \cdot |F_2| \cdot \deg(F_1) \cdot \deg(F_2) \cdot (\deg(F_1) + \deg(F_2))).$$

2.4. Classical graph similarity measures

In this section we do not give a complete review of graph similarity measuring, but we repeat only results relevant for this paper. First, we give a well-known key result due to Zelinka [25].

Theorem 2.1. Let $\mathcal{H}, \widetilde{\mathcal{H}}$ be unlabeled graphs without reflexive² and multiple edges. Furthermore let $|V| = |\widetilde{V}| = n$. $\overline{SUB}_m(\mathcal{H})$ denotes the set of induced subgraphs of order m . \mathcal{H}^\star denotes the isomorphism classes of such graphs in which lie \mathcal{H} and let

$$SUB_m(\mathcal{H}) := \{\mathcal{H}^\star | \mathcal{H} \in \overline{SUB}_m(\mathcal{H})\}. \quad (1)$$

$SUB_m(\mathcal{H})$ is just the set of isomorphism classes in which lie the induced subgraphs of \mathcal{H} with order m . Then,

$$d_Z(\mathcal{H}, \widetilde{\mathcal{H}}) := n - \text{SIM}(\mathcal{H}, \widetilde{\mathcal{H}}), \quad (2)$$

is a graph metric, where

$$\text{SIM}(\mathcal{H}, \widetilde{\mathcal{H}}) := \max\{m | SUB_m(\mathcal{H}) \cap SUB_m(\widetilde{\mathcal{H}}) \neq \emptyset\}. \quad (3)$$

Sobik [20] and Kaden [12] generalized this measure for arbitrary (also labeled) graphs of different orders.

Theorem 2.2. Let $\mathcal{H} := (V, E, f_V, f_E, A_V, A_E)$ be a finite, labeled and directed graph. A_V, A_E denote finite, non-empty vertex and edge alphabets and $f_V: V \rightarrow A_V, f_E: E \rightarrow A_E$ the associated vertex and edge labeling functions. Now, let \mathcal{H} and $\widetilde{\mathcal{H}}$ be finite, labeled graphs of arbitrary orders. Then,

$$d_S(\mathcal{H}, \widetilde{\mathcal{H}}) := \max\{|\mathcal{H}|, |\widetilde{\mathcal{H}}|\} - \text{SIM}(\mathcal{H}, \widetilde{\mathcal{H}}) \quad (4)$$

is a graph metric.

Fig. 2 shows the situation for determining subgraph isomorphism. Fig. 3 shows the application of Theorem 2.2. A direct application of the similarity measure of Zelinka provides Kaden [12]. Kaden transforms graphs in line graphs [1] and applies the theorem of Zelinka to the transformed graphs. For example the line graph of an undirected graph $G = (V, E)$ is defined by $\overline{G} = (E, \overline{E}), \overline{E} := \{e, \hat{e} | e, \hat{e} \in E \text{ and } e, \hat{e} \text{ are incident in } G\}$. By iterating the application of the line graph definition, Kaden obtains the graphs \overline{G}^m . As a result, it holds

Theorem 2.3. Let A_n be the set of connected graphs of order n . For given $0 \leq m < n$ it holds:

$$d_K^m(\mathcal{H}, \widetilde{\mathcal{H}}) = d_S(\mathcal{H}^m, \widetilde{\mathcal{H}}^m) = \max\{|\mathcal{H}^m|, |\widetilde{\mathcal{H}}^m|\} - \text{SIM}(\mathcal{H}^m, \widetilde{\mathcal{H}}^m), \quad (5)$$

is a graph metric on A_n .

Another well-known approach to measure the structural similarity of graphs is due to Bunke et al. [3].

Definition 2.6. We denote as an Optimal Inexact Match a sequence of transformations (insertions, deletions and substitutions of vertices) which transform \mathcal{H}_1 to \mathcal{H}_2 by producing minimal transformation costs. Assuming m_1, m_2, \dots, m_n are all possible Inexact Matches between \mathcal{H}_1 and \mathcal{H}_2 . Then, the Optimal Inexact Match m' is defined by $c(m') = \min\{c(m_i) | 1 \leq i \leq n\}$, where $c(m_i)$ denotes the costs of c_i .

Theorem 2.4. Let $d(\mathcal{H}_1, \mathcal{H}_2)$ be the costs for determining the Optimal Inexact Match between \mathcal{H}_1 and \mathcal{H}_2 . $d(\mathcal{H}_1, \mathcal{H}_2)$ is a graph metric.

A thorough overview of graph similarity measures can be found in [4,5].

2.5. Computational complexity of classical graph similarity measures

From the algorithm of Ullman [23] follows that the subgraph isomorphism problem is \mathcal{NP} -complete for general graphs. Then, we receive that the computation of Zelinka's measure and the generalizations of Sobik and Kaden are \mathcal{NP} -complete. The theorem of Bunke et al. leads us to the same result.

² That means edges of the type $(v, v) \in E$.

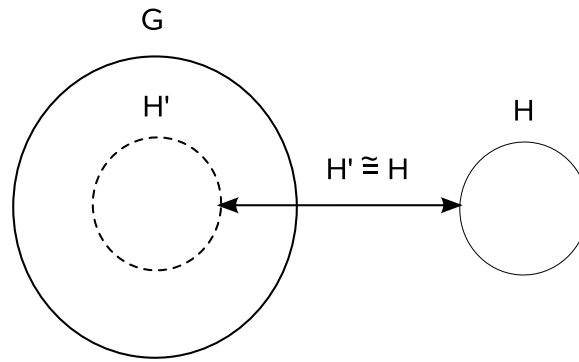


Fig. 2. The graphs $G, H, |G| > |H|$ and H' are illustrated as sets.

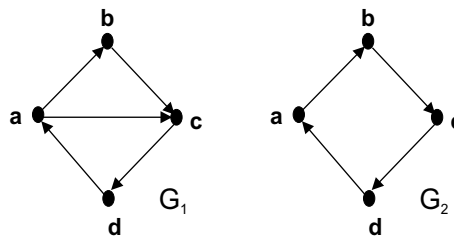


Fig. 3. Two vertex labeled directed graphs. The application of [Theorem 2.2](#) yields $d_S(G_1, G_2) = 4 - 3 = 1$.

3. Novel approach to measure the similarity of graphs

In this section we repeat our novel approach to measure the structural similarity of generalized trees [6]. Our novel approach is based on the following idea: Transform the underlying graphs into property strings and then align these strings with a dynamic programming [2] algorithm. Hence, we reduce a graph similarity problem into a string similarity problem. From the used string alignment technique we obtain similarity measures which can be computed polynomially in time.

For explaining the construction of our novel method we note that the graphs under consideration – generalized trees – can be transformed into property strings in a level oriented way. The property strings are the out-degree and in-degree sequences on each level of the generalized trees. In order to solve the graph similarity problem for our graph class we consider the transformations

$$\begin{aligned}
 S_0^{\mathcal{H}_1} &:= r_1^{\mathcal{H}_1}, \\
 S_1^{\mathcal{H}_1} &:= v_{1,1}^{\mathcal{H}_1} \circ v_{1,2}^{\mathcal{H}_1} \circ \dots \circ v_{1, \delta_{\text{out}}(r_1^{\mathcal{H}_1})}^{\mathcal{H}_1}, \\
 &\vdots \\
 S_{h_1}^{\mathcal{H}_1} &:= v_{h_1,1}^{\mathcal{H}_1} \circ v_{h_1,2}^{\mathcal{H}_1} \circ \dots \circ v_{h_1, \sigma_{h_1}}^{\mathcal{H}_1}, \\
 S_0^{\mathcal{H}_2} &:= r_2^{\mathcal{H}_2}, \\
 S_1^{\mathcal{H}_2} &:= v_{1,1}^{\mathcal{H}_2} \circ v_{1,2}^{\mathcal{H}_2} \circ \dots \circ v_{1, \delta_{\text{out}}(r_2^{\mathcal{H}_2})}^{\mathcal{H}_2}, \\
 &\vdots \\
 S_{h_2}^{\mathcal{H}_2} &:= v_{h_2,1}^{\mathcal{H}_2} \circ v_{h_2,2}^{\mathcal{H}_2} \circ \dots \circ v_{h_2, \sigma_{h_2}}^{\mathcal{H}_2},
 \end{aligned}$$

representing generalized trees $\mathcal{H}_1, \mathcal{H}_2$ together with their induced out-degree and in-degree sequences on a level i (see also [Definition 2.1](#)). It holds $r_k^{\mathcal{H}_k} := v_{0,1}^{\mathcal{H}_k}, k \in \{0, 1\}$. The problem of determining the structural similarity between \mathcal{H}_1 and \mathcal{H}_2 is now defined as alignment of the sequences above according to a cost function α .

In other words: the more similar the out-degree and in-degree alignments on a level i , $0 \leq i \leq h$ are, the more similar is the common structure of the graphs, and vice versa.

Fig. 4 shows two generalized trees together with their property strings on the graph levels. Now, we can write the sequences $S_0^{\mathcal{H}_1}, S_1^{\mathcal{H}_1}, \dots, S_k^{\mathcal{H}_1}, k \in \{1, 2\}$ in a more compact form

$$S_1 := r_1^{\mathcal{H}_1} \circ v_{1,1}^{\mathcal{H}_1} \circ v_{1,2}^{\mathcal{H}_1} \circ \dots \circ v_{h_1, \sigma_{h_1}}^{\mathcal{H}_1}, \tag{6}$$

$$S_2 := r_2^{\mathcal{H}_2} \circ v_{1,1}^{\mathcal{H}_2} \circ v_{1,2}^{\mathcal{H}_2} \circ \dots \circ v_{h_2, \sigma_{h_2}}^{\mathcal{H}_2}. \tag{7}$$

That means we transform the generalized trees $\mathcal{H}_1, \mathcal{H}_2$ in sequences S_1, S_2 which have not necessarily the same length. Here, we assume $S_k[i]$ as the i th position of the sequences S_k and it holds $S_1[n] = v_{h_1, \sigma_{h_1}}^{\mathcal{H}_1}, S_2[m] = v_{h_2, \sigma_{h_2}}^{\mathcal{H}_2}, \mathbb{N} \ni n, m \geq 1, S_k[1] = r_k^{\mathcal{H}_k}, k \in \{1, 2\}$.

For each possible aligned pair $[a, b]$ a cost function $\alpha([a, b]) \in \mathbb{R}_+$ is assigned, where a, b are sequence entries of S_1 and S_2 or the gap symbol $-$. The algorithm for finding the optimal alignment of the sequences (6) and (7) produces a matrix $(\mathcal{M}(i, j))_{ij}, 0 \leq i \leq n, 0 \leq j \leq m$, where $\mathcal{M}(i, j)$ is equivalent to the minimal edit distance between the sequences \tilde{S}_1, \tilde{S}_2 . \tilde{S}_1 and \tilde{S}_2 consist of the first i th, j th characters of S_1 and S_2 . Finally, we are looking for $\mathcal{M}(n, m)$ because $\mathcal{M}(n, m)$ is the minimal edit distance

$$d(S_1, S_2) := \min_{S_1 \rightarrow S_2} \sum \alpha([a, b]) \tag{8}$$

due to Levenstein [14]. This leads us to the following definition.

Definition 3.1. The recursive algorithm is defined by the following equations:

$$\mathcal{M}(0, 0) := 0, \tag{9}$$

$$\mathcal{M}(i, 0) := \mathcal{M}(i - 1, 0) + \alpha(S_1[i], -), \quad 1 \leq i \leq n, \tag{10}$$

$$\mathcal{M}(0, j) := \mathcal{M}(0, j - 1) + \alpha(-, S_2[j]), \quad 1 \leq j \leq m, \tag{11}$$

$$\mathcal{M}(i, j) := \min \begin{cases} \mathcal{M}(i - 1, j) + \alpha(S_1[i], -), \\ \mathcal{M}(i, j - 1) + \alpha(-, S_2[j]), \\ \mathcal{M}(i - 1, j - 1) + \alpha(S_1[i], S_2[j]), \end{cases} \quad i \in [1, n], \quad j \in [1, m]. \tag{12}$$

By tracing back along the minimal values from $\mathcal{M}(n, m)$ to $\mathcal{M}(0, 0)$ we find the optimal alignment between S_1 and S_2 . In Section 4 we will analyze the complexity of the graph similarity problem of generalized trees. Starting from this algorithm we construct now similarity functions $\alpha^{\text{out}}, \alpha^{\text{in}}$ in order to define the final graph similarity measures $(d_i(\mathcal{H}_1, \mathcal{H}_2))_{1 \leq i \leq k}$. We define:

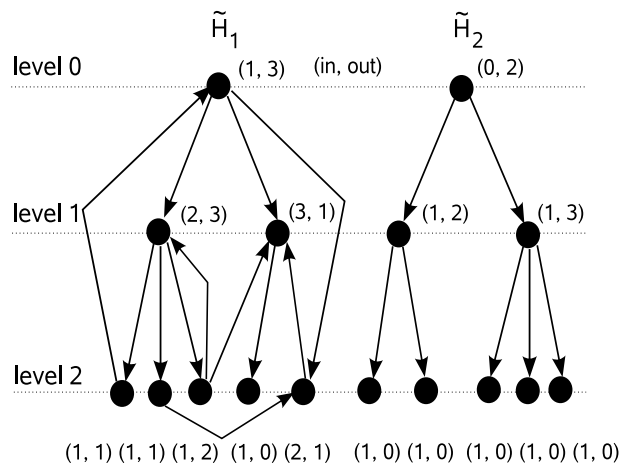


Fig. 4. Property strings in tuple-form: out-degree and in-degree sequences on the generalized tree levels.

$$\alpha^{\text{out}}(v_{i_1, j_1}^{\mathcal{H}_1}, v_{i_2, j_2}^{\mathcal{H}_2}) := \begin{cases} \beta^{\text{out}}(\delta_{\text{out}}(v_{i_1, j_1}^{\mathcal{H}_1}), \delta_{\text{out}}(v_{i_2, j_2}^{\mathcal{H}_2}), \sigma_{\text{out}}^1), & i_1 = i_2, \\ +\infty, & \text{else,} \end{cases} \quad (13)$$

$0 \leq i_k \leq h_k, 1 \leq j_k \leq \sigma_{i_k}^k$. Thereby, it holds $\beta^{\text{out}}(x, y, \sigma_{\text{out}}^k) := 1 - e^{-\frac{1}{2} \frac{(x-y)^2}{(\sigma_{\text{out}}^k)^2}}$, $x, y, \sigma_{\text{out}}^k \in \mathbb{R}, k \in \{1, 2\}$, and

$$\alpha^{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}, -) := \beta^{\text{out}}(\delta_{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}), \xi, \sigma_{\text{out}}^2), \quad (14)$$

$$\alpha^{\text{out}}(-, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{out}}(\xi, \delta_{\text{out}}(v_{i, j_2}^{\mathcal{H}_2}), \sigma_{\text{out}}^2). \quad (15)$$

By defining $\beta^{\text{in}}(x, y, \sigma_{\text{in}}^k) := 1 - e^{-\frac{1}{2} \frac{(x-y)^2}{(\sigma_{\text{in}}^k)^2}}$ we obtain in the same way

$$\alpha^{\text{in}}(v_{i_1, j_1}^{\mathcal{H}_1}, v_{i_2, j_2}^{\mathcal{H}_2}) := \begin{cases} \beta^{\text{in}}(\delta_{\text{in}}(v_{i_1, j_1}^{\mathcal{H}_1}), \delta_{\text{in}}(v_{i_2, j_2}^{\mathcal{H}_2}), \sigma_{\text{in}}^1), & i_1 = i_2, \\ +\infty, & \text{else,} \end{cases} \quad (16)$$

$$\alpha^{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}, -) := \beta^{\text{in}}(\delta_{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}), \xi, \sigma_{\text{in}}^2), \quad (17)$$

$$\alpha^{\text{in}}(-, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{in}}(\xi, \delta_{\text{in}}(v_{i, j_2}^{\mathcal{H}_2}), \sigma_{\text{in}}^2), \quad (18)$$

$\xi > 0$. The choice $\xi > 0$ in Eqs. (14), (15), (17) and (18) prevents an alignment between two leaves being better evaluated as an alignment between a leaf and a gap. The definitions (13) and (16) of the functions $\alpha^{\text{out}}, \alpha^{\text{in}}$ state that we do not align vertices on different levels. To prevent this, we set the gap penalty to $+\infty$, whereby our dynamic programming algorithm will never choose this cost-intensive path. In order to evaluate the alignments on each level of the given graphs $\mathcal{H}_1, \mathcal{H}_2$, we define the functions

$$\gamma^{\text{out}}(\mathcal{H}_k, i) := \frac{\sum_{j=1}^{\sigma_i^k} \hat{\alpha}^{\text{out}}(v_{i, j}^{\mathcal{H}_k}, \text{align}(v_{i, j}^{\mathcal{H}_k}))}{\sigma_i^k}, \quad (19)$$

$$\gamma^{\text{in}}(\mathcal{H}_k, i) := \frac{\sum_{j=1}^{\sigma_i^k} \hat{\alpha}^{\text{in}}(v_{i, j}^{\mathcal{H}_k}, \text{align}(v_{i, j}^{\mathcal{H}_k}))}{\sigma_i^k}, \quad (20)$$

$k \in \{1, 2\}$. Once again, σ_i^k denotes the upper index of a vertex on level i related to \mathcal{H}_k . For completing the definitions (19), (20), we define the mapping align and $\hat{\alpha}^{\text{out}}, \hat{\alpha}^{\text{in}}$ as follows:

$$\text{align}(v_{i, j_1}^{\mathcal{H}_1}) := \begin{cases} v_{i, j_2}^{\mathcal{H}_2}, & \text{align}^{-1}(v_{i, j_2}^{\mathcal{H}_2}) = v_{i, j_1}^{\mathcal{H}_1}, \\ -, & \text{else,} \end{cases} \quad (21)$$

$$\hat{\alpha}^{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}, -) := \beta^{\text{out}}(\delta_{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}), \xi, \hat{\sigma}_{\text{out}}^1), \quad (22)$$

$$\hat{\alpha}^{\text{out}}(-, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{out}}(\xi, \delta_{\text{out}}(v_{i, j_2}^{\mathcal{H}_2}), \hat{\sigma}_{\text{out}}^1), \quad (23)$$

$$\hat{\alpha}^{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{out}}(\delta_{\text{out}}(v_{i, j_1}^{\mathcal{H}_1}), \delta_{\text{out}}(v_{i, j_2}^{\mathcal{H}_2}), \hat{\sigma}_{\text{out}}^2), \quad (24)$$

$$\hat{\alpha}^{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}, -) := \beta^{\text{in}}(\delta_{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}), \xi, \hat{\sigma}_{\text{in}}^1), \quad (25)$$

$$\hat{\alpha}^{\text{in}}(-, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{in}}(\xi, \delta_{\text{in}}(v_{i, j_2}^{\mathcal{H}_2}), \hat{\sigma}_{\text{in}}^1), \quad (26)$$

$$\hat{\alpha}^{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}, v_{i, j_2}^{\mathcal{H}_2}) := \beta^{\text{in}}(\delta_{\text{in}}(v_{i, j_1}^{\mathcal{H}_1}), \delta_{\text{out}}(v_{i, j_2}^{\mathcal{H}_2}), \hat{\sigma}_{\text{in}}^2). \quad (27)$$

In order to define local similarity functions which weighs out-degree and in-degree sequences on a generalized tree level, we set

$$\gamma^{\text{out}}(i) := 1 - \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{out}}(v_{i, j}^{\mathcal{H}_1}, \text{align}(v_{i, j}^{\mathcal{H}_1})) \right\} + \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{out}}(v_{i, j}^{\mathcal{H}_2}, \text{align}(v_{i, j}^{\mathcal{H}_2})) \right\}, \quad (28)$$

and

$$\gamma^{\text{in}}(i) := 1 - \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{in}}\left(v_{i,j}^{\mathcal{H}_1}, \text{align}\left(v_{i,j}^{\mathcal{H}_1}\right)\right) \right\} + \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{in}}\left(v_{i,j}^{\mathcal{H}_2}, \text{align}\left(v_{i,j}^{\mathcal{H}_2}\right)\right) \right\}. \tag{29}$$

Hence, we obtain measures which indicate how similar the out-degree and in-degree alignments of two sequences on a level i are. Now, we get

Corollary 3.1

$$\gamma^{\text{out}}(i), \gamma^{\text{in}}(i) \in [0, 1].$$

In order to classify our final graph similarity measures we first express

Definition 3.2. Let U be a set of units and a mapping $\phi : U \times U \rightarrow [0, 1]$. ϕ is called a Backward similarity measure if

$$\phi(u, v) = \phi(v, u) \quad \forall u, v \in U \quad (\text{Symmetry}) \tag{30}$$

and

$$\phi(u, u) \geq \phi(u, v) \quad \forall u, v \in U \quad (\text{Backward}). \tag{31}$$

Finally, we prove the main result for measuring the structural similarity of generalized trees.

Theorem 3.2. Let $\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2$ be generalized trees, $0 \leq i \leq \rho$, $\rho := \max(h_1, h_2)$.

$$d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) := \frac{\sum_{i=0}^{\rho} \lambda_i \cdot \gamma^{\text{fin}}(i)}{\sum_{i=0}^{\rho} \lambda_i}, \tag{32}$$

$$d_2(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) := \frac{\sum_{i=0}^{\rho} \gamma^{\text{fin}}(i)}{\rho + 1}, \tag{33}$$

$$d_3(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) := \frac{\prod_{i=0}^{\rho} \gamma^{\text{fin}}(i)}{d_2(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2)}, \tag{34}$$

are a family $(d_i(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2))_{1 \leq i \leq 3}$ of Backward similarity measures, where $\gamma^{\text{fin}}(i)$ is defined as

$$\gamma^{\text{fin}}(i) := \zeta \cdot \gamma^{\text{out}}(i) + (1 - \zeta) \cdot \gamma^{\text{in}}(i). \tag{35}$$

Proposition 3.3. It holds $(d_i(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2))_{1 \leq i \leq 3} \in [0, 1]$.

Proof. Starting from Eqs. (32)–(34), Proposition 3.3 follows immediately from Corollary 3.1. \square

Proof of Theorem 3.2. To prove we first consider the function

$$\gamma^{\text{out}}(i) := 1 - \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{out}}\left(v_{i,j}^{\mathcal{H}_1}, \text{align}\left(v_{i,j}^{\mathcal{H}_1}\right)\right) \right\} + \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{out}}\left(v_{i,j}^{\mathcal{H}_2}, \text{align}\left(v_{i,j}^{\mathcal{H}_2}\right)\right) \right\}.$$

With Eq. (21), it follows that we have to distinguish the three cases for the function $\hat{\alpha}^{\text{out}} : \hat{\alpha}^{\text{out}}\left(v_{i,j_1}^{\mathcal{H}_1}, -\right), \hat{\alpha}^{\text{out}}\left(-, v_{i,j_2}^{\mathcal{H}_2}\right), \hat{\alpha}^{\text{out}}\left(v_{i,j_1}^{\mathcal{H}_1}, v_{i,j_2}^{\mathcal{H}_2}\right)$. We infer by Eqs. (22)–(24) that

$$\hat{\alpha}^{\text{out}}\left(v_{i,j_1}^{\mathcal{H}_1}, \text{align}\left(v_{i,j_1}^{\mathcal{H}_1}\right)\right) \leq 1 \quad \text{and} \quad \hat{\alpha}^{\text{out}}\left(v_{i,j_2}^{\mathcal{H}_2}, \text{align}\left(v_{i,j_2}^{\mathcal{H}_2}\right)\right) \leq 1.$$

Hence, by definition (28) we obtain $\gamma^{\text{out}}(i) \leq 1$. The proof $\gamma^{\text{in}}(i) \leq 1$ is identical. Since

$$\gamma^{\text{fin}}(i) \leq \zeta + (1 - \zeta) = 1,$$

we obtain

$$d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) \leq \frac{\sum_{i=0}^{\rho} \lambda_i}{\sum_{i=0}^{\rho} \lambda_i} = 1. \tag{36}$$

To prove that $d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2)$ is symmetric, we see that

$$\begin{aligned} \gamma^{\text{out}}(i) &:= 1 - \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{out}}(v_{i,j}^{\hat{\mathcal{H}}_1}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_1})) \right\} + \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{out}}(v_{i,j}^{\hat{\mathcal{H}}_2}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_2})) \right\} \\ &= 1 - \frac{1}{\sigma_i^2 + \sigma_i^1} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{out}}(v_{i,j}^{\hat{\mathcal{H}}_2}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_2})) \right\} + \frac{1}{\sigma_i^2 + \sigma_i^1} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{out}}(v_{i,j}^{\hat{\mathcal{H}}_1}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_1})) \right\}, \end{aligned}$$

and

$$\begin{aligned} \gamma^{\text{in}}(i) &:= 1 - \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{in}}(v_{i,j}^{\hat{\mathcal{H}}_1}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_1})) \right\} + \frac{1}{\sigma_i^1 + \sigma_i^2} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{in}}(v_{i,j}^{\hat{\mathcal{H}}_2}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_2})) \right\} \\ &= 1 - \frac{1}{\sigma_i^2 + \sigma_i^1} \cdot \left\{ \sum_{j=1}^{\sigma_i^1} \hat{\alpha}^{\text{in}}(v_{i,j}^{\hat{\mathcal{H}}_2}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_2})) \right\} + \frac{1}{\sigma_i^2 + \sigma_i^1} \cdot \left\{ \sum_{j=1}^{\sigma_i^2} \hat{\alpha}^{\text{in}}(v_{i,j}^{\hat{\mathcal{H}}_1}, \text{align}(v_{i,j}^{\hat{\mathcal{H}}_1})) \right\}. \end{aligned}$$

Therefore, we conclude from Eq. (32) that

$$d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2) = d_1(\hat{\mathcal{H}}_2, \hat{\mathcal{H}}_1).$$

For finalizing the proof for the similarity measure (32), we have to show that

$$d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_1) \geq d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2).$$

If $\hat{\mathcal{H}}_1 = \hat{\mathcal{H}}_2$ then, $\gamma^{\text{out}}(i) = 1$, $\gamma^{\text{in}}(i) = 1$ and $\gamma^{\text{fin}}(i) = 1$. Therefore we infer from Eq. (32) that $d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_1) = 1$ and see

$$d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_1) = 1 \geq \frac{\sum_{i=0}^{\rho} \lambda_i \cdot \gamma^{\text{fin}}(i)}{\sum_{i=0}^{\rho} \lambda_i} = d_1(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2),$$

the Backward property. In the case of similarity measure (33), we have nothing to prove, because if we set $1 = \lambda_0 = \lambda_1 = \dots = \lambda_{\rho}$ in Eq. (32), we obtain Eq. (33). To prove the assertion of the theorem for Eq. (34), we consider the well-known inequality

$$(p_1 \cdot p_2 \cdot \dots \cdot p_n)^{\frac{1}{n}} \leq \frac{p_1 + p_2 + \dots + p_n}{n}, \quad p_i > 0, \quad 1 \leq i \leq n. \tag{37}$$

Since $\gamma^{\text{fin}}(i) \leq 1$, we can apply inequality (37) and get

$$\gamma^{\text{fin}}(0) \cdot \gamma^{\text{fin}}(1) \cdot \dots \cdot \gamma^{\text{fin}}(\rho) \leq [\gamma^{\text{fin}}(0) \cdot \gamma^{\text{fin}}(1) \cdot \dots \cdot \gamma^{\text{fin}}(\rho)]^{\frac{1}{\rho+1}} \leq \frac{\gamma^{\text{fin}}(0) + \gamma^{\text{fin}}(1) + \dots + \gamma^{\text{fin}}(\rho)}{\rho + 1}.$$

Especially, we have

$$1 \geq (\rho + 1) \frac{\gamma^{\text{fin}}(0) \cdot \gamma^{\text{fin}}(1) \cdot \dots \cdot \gamma^{\text{fin}}(\rho)}{\gamma^{\text{fin}}(0) + \gamma^{\text{fin}}(1) + \dots + \gamma^{\text{fin}}(\rho)}. \tag{38}$$

The symmetry condition is clear, because the expression in the denominator of Eq. (34) is a special case of d_1 . The Backward condition follows immediately from Inequality Eq. (38),

$$1 = d_3(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_1) \geq \frac{\prod_{i=0}^{\rho} \gamma^{\text{fin}}(i)}{d_2(\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2)}. \quad \square$$

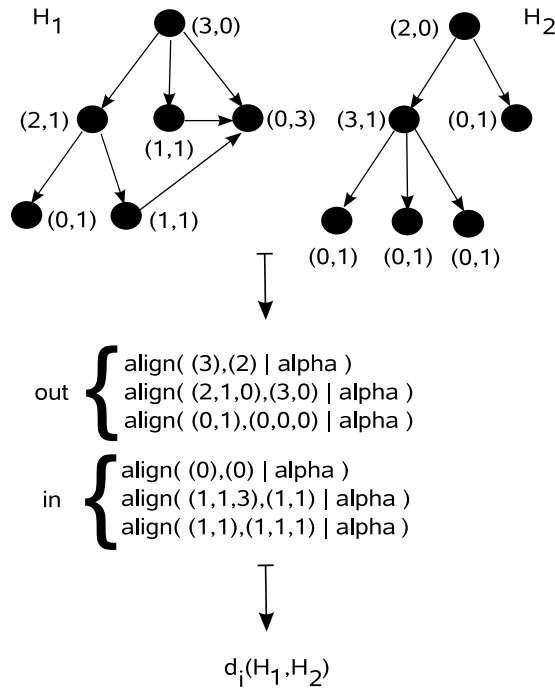


Fig. 5. Steps of similarity measuring for generalized trees.

Fig. 5 shows the overall approach exemplarily. In the following Section 4 we analyze the computational complexity of our novel graph similarity measure.

4. Efficiency analysis

To determine the computational complexity, first, we need the following definition.

Definition 4.1. Let $V_{S_1, S_2} := \{(i, j) \mid 0 \leq i \leq n, 0 \leq j \leq m\}$ be the vertex set and the edge types are $e_{\text{Del}} := (i - 1, j) \rightarrow (i, j)$, $e_{\text{Ins}} := (i, j - 1) \rightarrow (i, j)$, $e_{\text{Subst}} := (i - 1, j - 1) \rightarrow (i, j)$. $f_{E_{S_1, S_2}} : E_{S_1, S_2} \rightarrow \mathbb{R}_+$ denotes the edge labeling function. The edge set is defined by

$$E_{S_1, S_2} := \begin{cases} e_{\text{Del}}, & i \in [1, n], f_{E_{S_1, S_2}}(e_{\text{Del}}) = [S_1[i], -], \\ e_{\text{Ins}}, & j \in [1, m], f_{E_{S_1, S_2}}(e_{\text{Ins}}) = [-, S_2[j]], \\ e_{\text{Subst}}, & i \in [1, n], j \in [1, m], f_{E_{S_1, S_2}}(e_{\text{Subst}}) = [S_1[i], S_2[j]], \end{cases}$$

$G_{S_1, S_2} := (V_{S_1, S_2}, E_{S_1, S_2}, f_{E_{S_1, S_2}})$ is called the alignment graph of the sequences S_1 and S_2 representing generalized trees.

The edge types of these graphs reflect computational meanings in terms of transforming S_1 to S_2 .

Definition 4.2. $(i - 1, j) \rightarrow (i, j)$ equals the deletion of $S_1[i]$ in S_1 , $(i, j - 1) \rightarrow (i, j)$ equals the insertion of $S_2[j]$ in S_1 at the i th position, and $(i - 1, j - 1) \rightarrow (i, j)$ equals the substitution of $S_1[i]$ by $S_2[j]$.

Definition 4.3. Let A be an arbitrary alphabet and $w_1 \in A^\star$. Let I, D, R, M be a sequence of edit operations. I (insert) denotes the insertion, D (deletion) denotes the deletion, R (replace) denotes the substitution and M (match) denotes the correspondence. A transformation of w_1 into w_2 based on these edit operations is called edit transcript.

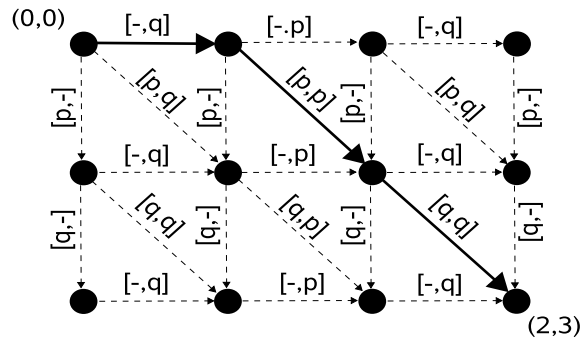


Fig. 6. The alignment graph G_{S_1, S_2} of S_1 and S_2 .

Definition 4.4. Let A be an arbitrary alphabet and $\{S_i | S_i \in A^{\star}\}$ be the set of finite sequences of A^{\star} . Then, $l : S \rightarrow \mathbb{N}$ denotes the function that determines the length of $S_i \in S$.

In order to state the time complexity of the recursive algorithm defined by Definition 3.1 we need a theorem of Gusfield [8] that states the relationship between alignments and paths of the associated alignment graph.

Theorem 4.1. Let $S_1, S_2 \in S$ be sequences and it holds $l(S_1) = n, l(S_2) = m$. The following assertions are equivalent:

- An edit transcript of S_1 and S_2 possesses a minimal number of edit operations.
- It corresponds a path in the associated alignment graph from vertex position $(0,0)$ to (n,m) with total minimal costs.

It follows straightforward

Corollary 4.2. The set of paths with total minimal costs from vertex position $(0,0)$ to (n,m) in the associated alignment graph specifies the set of optimal edit transcripts of S_1 and S_2 . The set of optimal edit transcripts describes the set of all optimal alignments from S_1 and S_2 .

Fig. 6 shows the alignment graph G_{S_1, S_2} of a simple example. If we set $S_1 = p \circ q$ and $S_2 = q \circ p \circ q$ the alignment is

$$\begin{array}{r} - \quad p \quad q \\ q \quad p \quad q \end{array}$$

Now, we obtain a theorem which states the computational complexity of the recursive algorithm defined by Definition 3.1.

Theorem 4.3. Let $\hat{\mathcal{H}}_1, \hat{\mathcal{H}}_2$ be generalized trees and S_1, S_2 be the associated vertex sequences, based on \hat{V}_1, \hat{V}_2 . The recursive algorithm for computing the optimal alignment between S_1 and S_2 given by Definition 3.1 possesses the computational complexity $\mathcal{O}(|\hat{V}_1| \cdot |\hat{V}_2| \cdot \mathcal{O}(1))$.

Proof. Starting from Definition 3.1 and especially with Eq. (12) we yield that the underlying algorithm needs time complexity $\mathcal{O}(|\hat{V}_1| \cdot |\hat{V}_2| \cdot \mathcal{O}(1))$ for creating the matrix $(\mathcal{M}(i,j))_{i,j}, 0 \leq i \leq |\hat{V}_1|, 0 \leq j \leq |\hat{V}_2|$. Furthermore, at each matrix position we compute the cost function α with time complexity $\mathcal{O}(1)$. It follows that the edit distance $\mathcal{M}(|\hat{V}_1|, |\hat{V}_2|)$ needs complexity $\mathcal{O}(|\hat{V}_1| \cdot |\hat{V}_2| \cdot \mathcal{O}(1))$. But Corollary 4.2 states that we have also found the optimal alignment. This completes the proof of the theorem. \square

Finally, we get

Corollary 4.4. *The structural similarity between two generalized trees $\hat{\mathcal{H}}_1$ and $\hat{\mathcal{H}}_2$ holds the time complexity $\mathcal{O}(|\hat{V}_1| \cdot |\hat{V}_2| \cdot \mathcal{O}(1))$.*

5. Conclusions

In the present paper we presented an algorithm with polynomial computational complexity to measure the structural similarity of a new graph class – generalized trees. The procedure to measure the structural similarity of graphs is completely different from known methods. First, we transformed the graphs in linear structures which contain structural information. Second, we determined the optimal alignment between the linear structures (out-degree and in-degree sequences on each level) and then we defined a family $(d_i)_{1 \leq i \leq 3}$ of graph similarity measures based on the obtained similarity scores (from $\gamma^{\text{out}}(i)$, $\gamma^{\text{in}}(i)$).

The efficiency analysis of the underlying algorithm was performed in Section 4. From this analysis, we found an efficient approach for graph similarity measuring in contrast to classical approaches [12,20,25] which are \mathcal{NP} -complete (see Section 2.5). Compared to known methods to measure the structural similarity of trees stated in Section 2.2 we found that our algorithm is comparable efficient despite its more general character.

References

- [1] J. Bang-Jensen, G. Gutin, Digraphs. Theory, Algorithms and Applications, Springer-Verlag, London–Berlin–Heidelberg, 2000.
- [2] R. Bellman, Dynamic Programming, Princeton University Press, 1957.
- [3] H. Bunke, G. Allermann, A metric on graphs for structural pattern recognition, in: H.W. Schussler (Ed.), SIGNAL PROCESSING II: Theory and Applications, 1983, pp. 257–260.
- [4] H. Bunke, Graph matching: theoretical foundations, algorithms, and applications, in: Proceedings of the Vision Interface 2000, Montreal/Canada, 2000, pp. 82–88.
- [5] H. Bunke, Recent developments in graph matching, in: Proceedings of the 15th International Conference on Pattern Recognition, Barcelona/Spain, 2000, vol. 2, pp. 117–124.
- [6] F. Emmert-Streib, M. Dehmer, J. Kilian, Classification of large graphs by a local tree decomposition, in: H.R. Arabnia, A. Scime (Eds.), Proceedings of the 2005 International Conference on Data Mining (DMIN'05), 2005, pp. 200–207.
- [7] R. Gleim, A. Mehler, M. Dehmer, Web corpus mining by instance of Wikipedia, Appears in Proceedings of the EACL 2006 Workshop on Web as Corpus, Trento, Italy, 3–7 April 2006.
- [8] D. Gusfield, Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology, Cambridge University Press, 1997.
- [9] M. Höchstmann, T. Töller, R. Giegerich, S. Kurtz, Local similarity in RNA secondary structures, in: Proceedings of the IEEE Computational Systems Bioinformatics Conference (CSB'03), 2003, pp. 159–168.
- [10] T. Jiang, L. Wang, K. Zhang, Alignment of Trees – An Alternative to Tree Edit, Theoretical Computer Science, vol. 143, Elsevier, 1995, pp. 137–148.
- [11] F. Kaden, Graph metrics and distance-graphs, in: M. Fiedler (Ed.), Graphs and other Combinatorial Topics Teubner Texte zur Math., Leipzig, 1983, vol. 59, pp. 145–158.
- [12] F. Kaden, Graphmetriken und Distanzgraphen, ZKI-Informationen, Akad. Wiss. DDR (1982) 1–63.
- [13] I. Koch, T. Lengauer, E. Wanke, An algorithm for finding maximal common subtopologies in a set of protein structures, J. Comput. Biol. 3 (2) (1996) 289–306.
- [14] V.I. Levenstein, Binary codes capable of correcting deletions insertions and reversals, Soviet Phys. Dokl. 10 (1966) 707–710.
- [15] A. Mehler, Hierarchical orderings of textual units, in: Proceedings of the 19th International Conference on Computational Linguistics (COLING 2002), Taipei, Taiwan, Morgan Kaufmann, San Francisco, 2002, pp. 646–652.
- [16] A. Mehler, M. Dehmer, R. Gleim, Towards logical hypertext structure. A graph-theoretic perspective, in: Proceedings of I2CS'04, Guadalajara/Mexico, Lecture Notes in Computer Science, Springer, Berlin–New York, 2004.
- [17] D. Sankoff, J.B. Kruskal, Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison, Addison-Wesley, 1983.
- [18] D. Sankoff, J.B. Kruskal, S. Mainville, R.J. Cedergren, Fast algorithms to determine RNA secondary structures containing multiple loops, in: D. Sankoff, J.B. Kruskal (Eds.), Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison, Addison-Wesley, 1983, pp. 93–120.
- [19] S.M. Selkow, The tree-to-tree editing problem, Inform. Process. Lett. 6 (6) (1977) 184–186.
- [20] F. Sobik, Graphmetriken und klassifikation strukturierter objekte, ZKI-Informationen Akad. Wiss. DDR 2 (82) (1982) 63–122.
- [21] F. Sobik, Modellierung von Vergleichsprozessen auf der Grundlage von Ähnlichkeitsmaßen für Graphen, ZKI-Informationen, Akad. Wiss. DDR 4 (1986) 104–144.

- [22] K. Tai, The tree-to-tree correction problem, *J. ACM* 26 (1979) 422–433.
- [23] J.R. Ullman, An algorithm for subgraph isomorphism, *J. ACM* 23 (1) (1976) 31–42.
- [24] L. Wang, J. Zhao, Parametric alignments of ordered trees, *Bioinformatics* 19 (17) (2003) 2237–2245.
- [25] B. Zelinka, On a certain distance between isomorphism classes of graphs, *Časopis pro řest. Matematiky* 100 (1975) 371–373.
- [26] K. Zhang, R. Statman, D. Shasha, On the editing distance between unordered labeled trees, *Int. Inform. Process. Lett.* 42 (3) (1992) 133–139.